

مقدمه

۱-۱- مقدمه

امروزه پیشرفت در تکنولوژی کامپیوترها، تقریباً همه جنبه‌های جامعه را تحت تأثیر قرار داده است. پیشرفت‌های صورت گرفته در سخت‌افزار، به برنامه‌نویسان اجازه داده است که برنامه‌های واقعاً مفیدی ایجاد کنند. برنامه‌نویسان موفق همیشه به سرعت اجرای برنامه‌های خود توجه می‌کنند چون دادن نتیجه سریع به کاربر در موفقیت یک برنامه نقش حیاتی دارد. در دهه‌های ۱۹۶۰ و ۱۹۷۰ اصلی‌ترین محدودیت در سرعت کامپیوترها اندازه حافظه کامپیوترها بود. بعدها پیشرفت‌های ایجاد شده در زمینه طراحی کامپیوترها و تکنولوژی حافظه‌ها واقعاً اهمیت مسأله حافظه کم را کاهش داد. امروزه برنامه‌نویسانی که به سرعت اجرای برنامه‌های خود علاقه دارند، باید به مسائلی توجه کنند که جایگزین مسأله کمبود حافظه شده‌اند. از جمله این مسائل می‌توان به ساختار سلسله مراتبی حافظه‌ها و موازات در پردازنده‌ها اشاره نمود.

برنامه‌نویسانی که در جستجوی راههایی هستند که کامپایلرها، سیستم عاملها، برنامه‌های پایگاه داده و حتی application هایی بسازند که قابل رقابت با محصول دیگران باشد، باید دانش خود را در ارتباط با ساختار و سازمان کامپیوتر افزایش دهند. ما در این کتاب توضیح خواهیم داد که داخل یک کامپیوتر چه چیزهایی وجود دارد و شما را با رموز برنامه‌نویسی آشنا خواهیم ساخت. همچنین توضیح داده خواهد شد که ساختار داخلی یک کامپیوتر چگونه بر کارایی برنامه‌ها تأثیر می‌گذارد. مهمتر از همه اینکه شما یاد خواهید گرفت که چگونه یک کامپیوتر برای خود بسازید.

فصل اول زیر بنای این کتاب است. این بخش به مطالب اساسی و تعاریف می‌پردازد. همچنین بخش‌های اصلی نرم‌افزار و سخت‌افزار را بیان می‌نماید. در این بخش مقدمه‌ای در مورد مدارهای فشرده (Integrated circuit) که تکنولوژی مؤثر در پیشرفت کامپیوترها بوده، خواهد آمد.

۱-۲- لایه‌های زیرین برنامه

برای اینکه واقعاً بتوانید با یک ماشین الکترونیکی مانند کامپیوتر ارتباط برقرار کنید باید از علامتهای الکترونیکی استفاده کنید. ساده‌ترین علائم قابل فهم برای این ماشین الکترونیکی، روشن (on) و خاموش (off) می‌باشد و بنابراین الفبای کامپیوتر دارای دو حرف می‌باشد. همان‌طور که ۲۶ حرف موجود در الفبای زبان انگلیسی محدودیتی بر مقدار مطالبی که نوشته می‌شود نمی‌گذارد، الفبای دو حرفی کامپیوتر نیز کارهایی که یک کامپیوتر انجام می‌دهد را محدود نمی‌کند. دو سمبلی که برای این دو حروف استفاده می‌شود ۰ و ۱ است و ما معمولاً زبان کامپیوتر را به صورت عددهایی در مبنای ۲ یا اعداد دودویی می‌شناسیم. به هر کدام از حروف (۰ و ۱) کامپیوتر یک رقم باینری یا بیت گفته می‌شود.

کامپیوترها فقط با فرمانهای ما کار می‌کنند، اسم هر فرمان تکی دستور (instruction) نامیده می‌شود. دستورات که مجموعه‌ای از بیت‌ها (۰ و ۱) می‌باشند که یک کامپیوتر می‌فهمد، می‌توانند به صورت اعداد در نظر گرفته شوند. به طور مثال مجموعه بیت زیر را در نظر بگیرید:

۱۰۰۰۱۱۰۰۱۰۱۰۰۰۰۰

این مجموعه بیت یک دستور می‌باشد که می‌تواند به یک کامپیوتر بگوید که دو عدد را با هم جمع کند. در فصل ۳ توضیح داده خواهد شد که چرا ما برای دستورات و داده‌ها از اعداد استفاده می‌کنیم. استفاده از اعداد برای دستورات و داده‌ها پایه و اساس محاسبات است بنابراین ما فصل ۳ را از دست نخواهیم داد!

اولین برنامه‌نویس‌های کامپیوتر، از طریق اعداد باینری با آن ارتباط برقرار می‌کردند یعنی به زبان ماشین کدنویسی می‌کردند. برنامه نویسی به زبان ماشین کار سختی بود و به همین دلیل به سرعت یک مجموعه علامتها اختراع شدند که به روش فکر کردن انسانها نزدیک بودند و برنامه نویسان از طریق این مجموعه علامتها برنامه نوشتند. اسمی که برای این زبان سمبلیک استفاده می‌شد، امروزه نیز کاربرد دارد و آن زبان اسمبلی است. در ابتدا برنامه نویسان کد نوشته شده با این علائم را به طور دستی به مجموعه ۰ و ۱ ها (زبان ماشین) ترجمه می‌کردند. ترجمه دستی نیز کار مشکلی بود. راهکار بعدی این بود که از خود کامپیوتر برای عمل ترجمه کمک بگیریم. برنامه‌نویسان اولیه برنامه‌هایی را جهت ترجمه از علامتهای نمادین به مجموعه ۰ و ۱ ها نوشتند و این برنامه‌ها را اسمبلر (assembler) نامگذاری کردند. اسمبلرها برنامه‌هایی بودند که یک برنامه را که به زبان اسمبلی نوشته شده بود در ورودی خود دریافت می‌کردند و در خروجی خود مجموعه‌ای از ۰ و ۱ ها تولید می‌کردند.

به طور مثال برنامه‌نویس می‌تواند دستوری را به صورت زیر بنویسد:

Add A, B

و اسمبلر می‌تواند این دستور را به صورت زیر ترجمه نماید:

۱۰۰۰۱۱۰۰۱۰۱۰۰۰۰۰

این دستور به کامپیوتر می‌گوید که دو عدد A و B را با هم جمع کند. با وجود ایجاد زبان اسمبلی، برنامه‌نویسی هنوز هم کار سختی بود. در زبان اسمبلی برنامه‌نویس باید برای هر دستور زبان ماشین که کامپیوتر می‌تواند آن را انجام دهد یک دستور اسمبلی بنویسد و این باعث می‌شود که برنامه‌نویس مانند یک ماشین فکر کند.

در اینجا یک سؤال ساده مطرح می‌شود:

اگر ما توانستیم یک برنامه‌ای بنویسیم که عمل ترجمه از زبان اسمبلی به زبان ماشین را انجام دهد، چه عاملی باعث می‌شود که نتوانیم برنامه‌ای بنویسیم که عمل ترجمه از نمادهای سطح بالاتر به زبان اسمبلی را انجام دهد.

جواب این است که هیچ چیز.

برنامه‌هایی که این نمادهای سطح بالاتر را قبول می‌کنند، کامپایلر نامیده می‌شوند و زبانهایی که آنها کامپایل می‌کنند، زبانهای برنامه‌نویسی سطح بالا نامیده می‌شوند. کامپایلرها یک برنامه‌نویس را قادر می‌سازند که عبارتی مانند عبارت زبان سطح بالای زیر را در برنامه‌های خود بنویسند.

$A+B$

یک کامپایلر ممکن است عبارت فوق را به عبارت زبان اسمبلی زیر کامپایل (ترجمه) نماید.

`Add A, B`

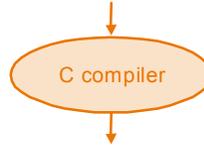
و یک اسمبلر ممکن است عبارت `add A, B` را به یک دستور باینری به صورت زیر ترجمه نماید که این دستور به کامپایلر بگوید دو عدد A و B را با هم جمع نماید.

۱۰۰۰۱۱۰۰۱۰۱۰۰۰۰۰

شکل ۱-۱ ارتباط بین این برنامه‌ها و زبانها را نشان می‌دهد.

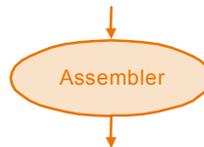
High-level
language
program
(in C)

```
swap(int v[], int k)
{int temp;
  temp = v[k];
  v[k] = v[k+1];
  v[k+1] = temp;
}
```



Assembly
language
program
(for MIPS)

```
swap:
  muli $2, $5, 4
  add $2, $4, $2
  lw $15, 0($2)
  lw $16, 4($2)
  sw $16, 0($2)
  sw $15, 4($2)
  jr $31
```



Binary machine
language
program
(for MIPS)

```
000000001010000100000000000011000
00000000100011100001100000100001
10001100011000100000000000000000
100011001111001000000000000000100
10101100111100100000000000000000
101011000110001000000000000000100
00000011111000000000000000001000
```

شکل ۱-۱ ارتباط بین زبانهای برنامه‌نویسی و مبدل‌های آنها.

زبانهای برنامه‌نویسی سطح بالا چندین مزیت مهم دارند. اول اینکه آنها به برنامه‌نویس قابلیت این را می‌دهند که در یک زبان خیلی نزدیک به زبان طبیعی فکر کنند. در این زبانها معمولاً برنامه‌نویس از کلمات و علامت‌های ریاضی زبان انگلیسی استفاده می‌کند. بعلاوه آنها اجازه می‌دهند زبان‌هایی متناظر با کاربرد مورد نظر طراحی شوند. به طور مثال فرترن (Fortran) برای محاسبات ریاضی، کوبول (Cobol) برای پردازش داده‌ها تجاری، لیسپ (Lisp) برای دستکاری سمبل‌ها و ... طراحی شده‌اند. دومین مزیت زبانهای برنامه‌نویسی سطح بالا افزایش تولیدات برنامه‌نویس است. در زبان‌های سطح بالا حجم کد کم است و تعداد زیادی کار را می‌توان در تعداد کمتری از خط‌های برنامه انجام داد این امر مزیت بزرگی نسبت به زبان اسمبلی که در آن حجم برنامه بسیار بزرگ می‌شود، می‌باشد.

آخرین مزیت این است که زبانهای برنامه‌سازی اجازه می‌دهند که برنامه‌ها از کامپیوتری که بر روی آن طراحی می‌شوند مستقل باشند. این به این دلیل است که کامپایلرها و اسمبلرها می‌توانند برنامه نوشته شده را به زبان ماشین کامپیوتری که می‌خواهد آن را اجرا کند تبدیل می‌کنند.

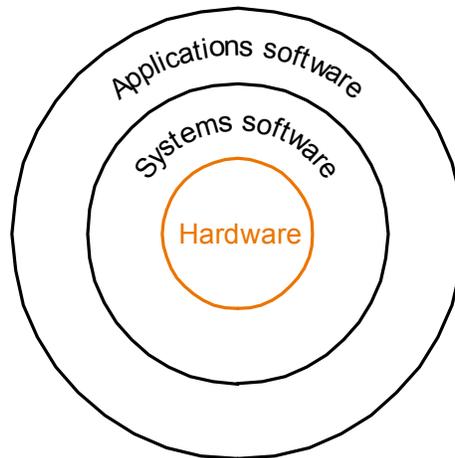
این سه مزیت آنقدر اهمیت دارند که امروزه تعداد برنامه‌های بسیار کمی با زبان اسمبلی نوشته می‌شود. با پیشرفت برنامه‌نویسی، برنامه‌نویسان مشاهده کردند که استفاده مجدد از برنامه‌ها، از اینکه همه چیز را از صفر شروع کنیم خیلی بهتر است. بنابراین برنامه‌نویسان این کار را شروع کردند که روتین‌های پرستفاده را در داخل کتابخانه‌ها (Libraries) جمع کنند. یکی از اولین کتابخانه‌های روتین‌ها، برای خواندن و نوشتن داده بود که به طور مثال شامل روتین‌هایی نظیر روتین‌های کنترل پرینترها بودند که به طور مثال موجود بودن کاغذ داخل پرینتر را قبل از عملیات پرینت کردن چک می‌کنند.

به زودی معلوم گردید که اگر یک برنامه‌ای وجود داشته باشد که اجرا شدن بقیه برنامه‌ها را مدیریت کند یک مجموعه از برنامه‌ها می‌توانند به طور مؤثری اجرا شوند، به محض اینکه اجرای یک برنامه به پایان رسید، برنامه‌ناظر، برنامه دیگری را از صف برنامه‌های منتظر اجرا، برای اجرا انتخاب می‌کند و بنابراین مانع ایجاد تأخیر می‌شود.

این برنامه ناظر که به زودی کتابخانه روتین‌های ورودی/خروجی نیز به آن اضافه شد، پایه و اساس چیزی بود که امروزه به آن سیستم عامل می‌گوییم. سیستم‌های عامل برنامه‌هایی هستند که منابع یک کامپیوتر را مدیریت می‌کنند تا اینکه برنامه‌های دیگر به طور مؤثری بر روی آن کامپیوتر اجرا شوند.

بعدها نرم‌افزارها بر اساس استفاده آنها دسته‌بندی شدند. نرم‌افزارهایی که معمولاً سرویس‌های مفید در اختیار کاربران قرار می‌دهند، نرم‌افزارهای سیستم (System software) نامیده می‌شوند. سیستم‌عامل‌ها، کامپایلرها و اسمبلرها مثالهایی از نرم‌افزارهای سیستم می‌باشند. در مقابل نرم‌افزارهای سیستمی، نرم‌افزارهای کاربردی (applications software) یا همان برنامه‌های کاربردی (application) قرار دارند که نامی است برای برنامه‌هایی که هدفشان کمک کردن به استفاده‌کننده‌های کامپیوتر است. به طور مثال از این دسته برنامه‌ها می‌توانیم به ادیتورهای متن اشاره کنیم.

شکل ۱-۲ لایه‌های سلسله مراتبی نرم‌افزار و جایگاه آنها نسبت به سخت‌افزار را نشان می‌دهد.



شکل ۱-۲: لایه‌های سلسله مراتبی نرم‌افزار و ارتباط آن با سخت‌افزار.

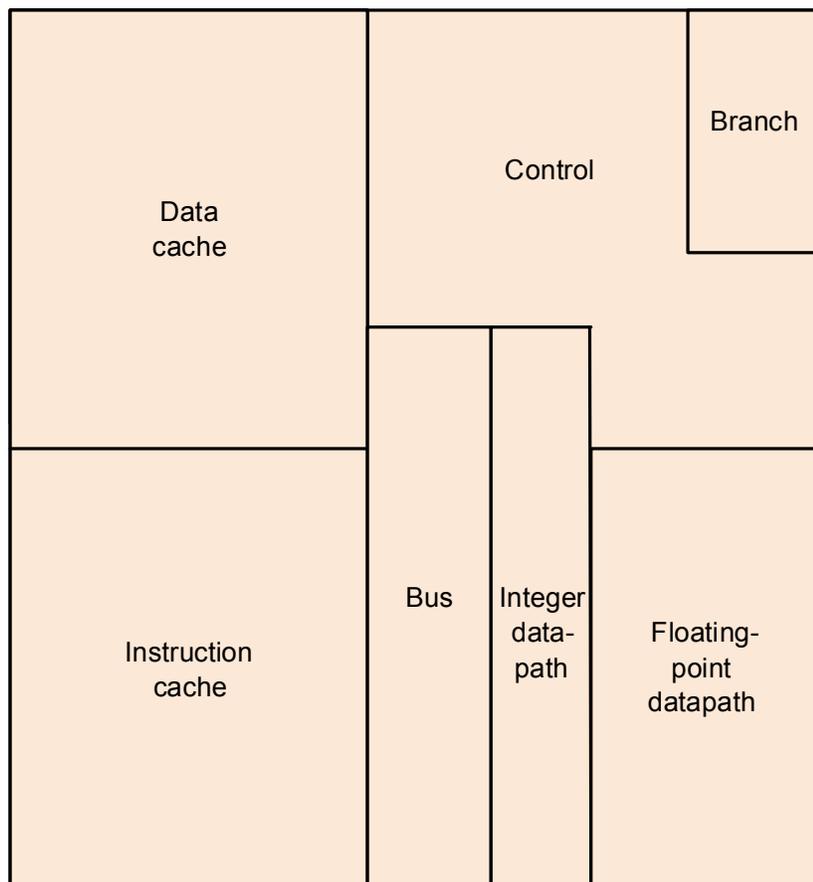
۱-۳- لایه‌های زیرین سخت‌افزار

در بخش قبل ما برای فهمیدن قسمت‌های زیرین یک نرم‌افزار به آنچه که در زیر برنامه شما قرار دارد نگاه انداختیم، حال می‌خواهیم چنین کاری را در رابطه با سخت‌افزار انجام دهیم. اگر یک کامپیوتر رومیزی را در نظر بگیریم می‌بینیم که دارای کیبورد، موس، صفحه نمایش و یک کیس که سخت‌افزارهای زیادی را در داخل خود نگه‌داری می‌کند، می‌باشد. این کامپیوتر معمولاً دارای یک اتصال شبکه نیز هست که ارتباط کامپیوتر را با کامپیوترهای دیگر، پرینتر و دیسکها برقرار می‌کند. دو ماجول اصلی کامپیوترها قطعات ورودی نظیر کیبورد و قطعات خروجی نظیر صفحه نمایش و پرینترها می‌باشند. همان طور که از نامشان پیدا است، با استفاده از ورودی، کامپیوتر تغذیه می‌شود و داده‌ها در اختیار آن قرار می‌گیرد و با استفاده از خروجی‌ها نتایج حاصل از محاسبات کامپیوتر از آن خارج شده و در اختیار کاربر قرار می‌گیرد. در فصل ۸ وسایل ورودی و خروجی با جزئیات بیشتری شرح داده خواهند شد.

اجازه دهید تا یک مرور مقدماتی به سخت‌افزار کامپیوتر داشته باشیم. وقتی که ما کیس یک کامپیوتر را باز می‌کنیم، یک بردی (board) را مشاهده می‌کنیم که با پلاستیک نازک سبز رنگی پوشیده شده است و بر روی آن تعداد زیادی مستطیل‌های کوچک سیاه و خاکستری قرار دارند. این برد، مادربرد (mother board) نام دارد. مستطیل‌های کوچک روی مادربرد IC (Integrated circuit) یا چیپ نامیده می‌شوند. این برد دارای سه قسمت مهم است، قسمتی که به قطعات ID (مثل درایو فلاپی و دیسک سخت) وصل می‌شود، حافظه و پردازنده. تعدادی برد دیگر نیز به مادربرد وصل می‌شوند که از آن جمله می‌توان به کارت شبکه و کارت گرافیکی اشاره نمود.

حافظه محلی است که برنامه‌ها به هنگام اجرا بر روی آن قرار می‌گیرند. همچنین حافظه داده‌های مورد نیاز برنامه در حال اجرا را نیز نگهداری می‌کند. حافظه معمولاً بر روی دو بورد کوچک که تقریباً در وسط مادربورد نصب می‌شوند، قرار می‌گیرد.

پردازنده قسمت فعال بورد است که دستورات برنامه را اجرا می‌کند. پردازنده اعداد را جمع می‌کند، تست می‌کند، به قطعات IO علامت می‌دهد تا فعال شوند و ... پردازنده به صورت یک مربع بزرگ بر روی مادربورد قرار دارد. معمولاً پردازنده را به اسم CPU می‌شناسیم که همان واحد پردازش مرکزی Central Processing Unit می‌باشد. شکل ۱-۳ محتویات داخلی یک CPU را نشان می‌دهد.



شکل ۱-۳: محتویات داخلی CPU

پردازنده از دو بخش اصلی تشکیل می‌شود:

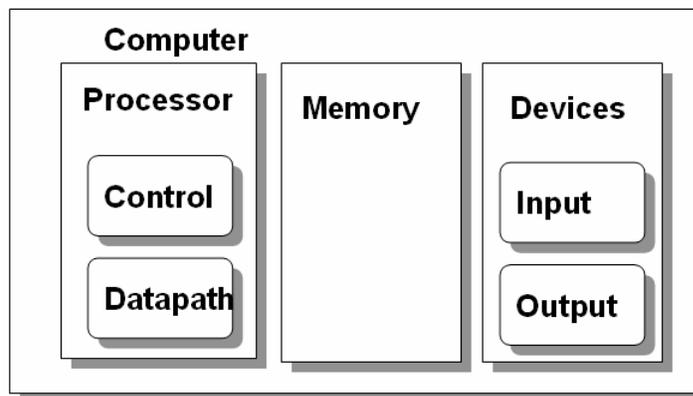
بخش *data path* و بخش کنترل (*Control*) که این دو قدرت و تفکر کامپیوتر را می‌سازند. *data path* عملیات ریاضی را انجام می‌دهد و بخش کنترل به *data path*، حافظه و قطعات IO دستور می‌دهد که

چه کاری انجام دهند. بر اساس اینکه دستوری که می‌خواهد انجام شود، چه عملیاتی می‌خواهد انجام دهد، واحد کنترل وظیفه سایر بخش‌ها را مشخص می‌کند.

فصل ۵ در مورد بخش‌های data path و کنترل صحبت خواهد کرد و فصل ۶ در مورد تغییراتی که باید در طراحی این دو بخش صورت بگیرد تا به سرعت بالاتری برسیم صحبت خواهد نمود.

شکل ۱-۳ مربوط به پردازنده پنتیوم اینتل می‌باشد. حجم سیلیکون استفاده شده برای این پردازنده ۹۱ میلی‌متر مربع بوده و در ساخت آن حدود ۳/۳ میلیون ترانزیستور استفاده شده است. حافظه کش cache این پردازنده تقریباً یک میلیون عدد از ۳ میلیون ترانزیستور را به خود اختصاص داده است. فصل ۷ توضیح خواهد داد که چرا کش منابع زیادی را مصرف می‌کند. قسمتهای دیگر چیپ در فصل-های بعدی توضیح داده خواهد شد. مدار پیش بینی دستور انشعاب (branch prediction) در فصل ۶ و گذرگاه سیستم Bus در فصل ۸ پوشش داده می‌شوند.

نکته مهم: ۵ بخش کلاسیک یک کامپیوتر عبارتند از ورودی، خروجی، حافظه data path و بخش کنترل که این دو بخش آخری معمولاً ترکیب شده و پردازنده (processor) نامیده می‌شود. شما می‌توانید هر بخشی از هر کامپیوتری را (چه کامپیوترهای قدیمی و چه کامپیوترهای فعلی) در یکی از این دسته‌ها قرار دهید. این پنج بخش در شکل ۱-۴ نشان داده شده‌اند.



شکل ۱-۴: پنج بخش کلاسیک یک کامپیوتر

پردازنده دستورات و داده‌ها را از حافظه دریافت می‌نماید؛ ورودی داده‌ها را در داخل حافظه می‌نویسد و خروجی داده‌ها را از حافظه می‌خواند. بخش کنترل، سیگنالهایی که عملیات data path، حافظه، ورودی و خروجی را مشخص می‌کنند، را تولید می‌نماید.

بر روی مادبورد معمولاً حافظه‌های نوع DRAM قرار دارد. DRAM علامت اختصاری عبارت Dynamic Random Access Memory می‌باشد. برای نگهداری داده‌ها و دستورات یک برنامه

چندین DRAM در کنار یکدیگر استفاده می‌شوند. در مقایسه با حافظه‌های با قابلیت دسترسی ترتیبی مانند نوارهای مغناطیسی که زمان دستیابی به قسمت‌های مختلف با هم فرق دارد، قسمت RAM از عبارت DRAM بیانگر این مطلب است که دسترسی به قسمت‌های مختلف این حافظه به یک اندازه طول می‌کشد. حافظه DRAM مانند یک بافر بزرگ یا یک آرایه بزرگ برای ذخیره داده‌ها مورد استفاده قرار می‌گیرد.

۱-۳-۱- معماری مجموعه دستورات

نکته‌ای که در اینجا باید به آن اشاره شود این است که ما هر چقدر وارد جزئیات سخت‌افزار و نرم‌افزار نشویم، به مدل‌های سطح بالای ساده‌ای از سخت‌افزار و نرم‌افزار خواهیم رسید. استفاده از این روش (مدل‌های لایه‌ای) یک تکنیک اساسی برای طراحی سیستم‌های کامپیوتری بسیار هوشمند می‌باشد. یکی از مهمترین سطوح تجرد^۱، واسط بین سخت‌افزار و نرم‌افزار سطح پایین^۲ می‌باشد. به دلیل اهمیت این موضوع یک اسم ویژه به آن اختصاص داده شده است. معماری مجموعه دستورات (ISA^۳)، یا به طور ساده معماری یک ماشین. ISA شامل همه چیزهایی که برنامه‌نویسان برای نوشتن برنامه‌های باینری زبان ماشین به آنها نیاز دارند می‌باشد. این نیازمندی‌ها شامل دستورات، قطعات I/O و غیره می‌باشد. اجزای مختلف یک معماری در فصل‌های ۳، ۴، ۷ و ۸ بحث خواهد شد. این واسط استاندارد به طراحان کامپیوتر اجازه می‌دهد تا مستقل از سخت‌افزار در مورد عملکردهای یک کامپیوتر صحبت کنند. به طور مثال ما می‌توانیم در مورد عملکرد یک ساعت دیجیتال (مانند نمایش زمان، زنگ اختار و ...) مستقل از اینکه سخت‌افزار آن چیست (صفحه نمایش، دکمه‌های پلاستیکی، چرخ دنده‌ها و ...) صحبت کنیم. به همین دلیل طراحان کامپیوتر بین یک پیاده‌سازی^۴ و یک معماری تفاوت قائل می‌شوند. یک پیاده‌سازی سخت‌افزاری است که سطوح تجرد معماری را فراهم می‌کند.

نکته مهم: نرم‌افزار و سخت‌افزار هر دو شامل لایه‌های سلسله مراتبی هستند، که لایه‌های زیرین جزئیات را از لایه‌های بالاتر مخفی نگه می‌دارند. پایه و اساس سطح تجرد این است که طراحان سخت‌افزار و نرم‌افزار پیچیدگی طراحی سیستم‌های کامپیوتری را کاهش دهند. یکی از واسط‌های کلیدی بین سطوح تجرد ISA می‌باشد: واسط بین سخت‌افزار و نرم‌افزار سطح پایین. این واسط امکان

1 - Abstract
2 - Low level software
Instruction Set Architecture-3
4 - Implementation

ایجاد پیاده‌سازی‌های مختلف با قیمت و سرعت متفاوت که نرم‌افزار یکسانی را اجرا می‌کنند در اختیار قرار می‌دهد.

استفاده‌کننده‌های macintosh اثر تغییر ISA را درک می‌کنند. برنامه‌هایی که برای معماری Power PC طراحی شده‌اند بر روی ماشین‌هایی که بر اساس 6800 ساخته شده‌اند اجرا نمی‌شوند، و برنامه‌هایی مبتنی بر 6800، به درستی بر روی Power PC اجرا نمی‌شوند. در مقابل خانواده 86*80 شرکت اینتل پیاده‌سازی‌های مختلفی از یک معماری را در اختیار قرار می‌دهند. برنامه‌های نوشته شده برای 8086 اولیه در سال ۱۹۷۸ می‌توانند بر روی Pentium Pro اجرا شوند. Intel در طی سالیان قابلیت‌های زیادی را به سیستم‌های قبلی افزوده است. اما همهٔ پردازنده‌های نسل بعدی برنامه‌های نوشته شده بر روی پردازنده‌های نسل‌های قبلی را اجرا می‌کنند.

۱-۳-۲- محل مطمئن برای نگهداری داده‌ها

اگر تغذیه (برق) کامپیوتر قطع شود، همه چیز از بین می‌رود و این به دلیل این است که حافظه‌های داخل کامپیوتر فرار هستند و با قطع برق اطلاعات آنها از بین می‌رود. در مقام مقایسه یک نوار کاست با قطع برق اطلاعات خود را از دست نمی‌دهد چون تکنولوژی آنها با حافظه کامپیوتر فرق می‌کند. در نوار کاست از قطعات مغناطیسی استفاده می‌شود.

کامپیوتر معمولاً دارای دو نوع حافظه است: حافظه اصلی و حافظه ثانوی

حافظه اصلی یا اولیه معمولاً از نوع DRAM بوده و اطلاعات آنها با قطع شدن برق از بین می‌رود ولی حافظه جانبی یا ثانویه از نوع مغناطیسی بوده و اطلاعات آنها با قطع شدن برق از بین نمی‌رود. بنابراین برای نگهداری داده‌ای که نباید از بین بروند از حافظه ثانویه استفاده می‌شود. به دلیل اینکه در حافظه‌های ثانویه از قطعات مکانیکی استفاده می‌شود، سرعت آنها نسبت به DRAM ها خیلی پایین‌تر است. زمان دسترسی به داده در حافظه ثانویه (هارد دیسک) معمولاً بین ۵ تا ۲۰ میلی ثانیه است در حالی که این زمان برای DRAM در حدود ۵۰ تا ۱۰۰ نانو ثانیه است. بنابراین DRAM حدود ۱۰۰۰۰۰ مرتبه سریعتر از هارد دیسک است.

هزینه تولید هارد دیسک‌ها نسبت به DRAM پایین‌تر است بنابراین قیمت آنها نسبت به DRAM کمتر است. در سال ۱۹۹۷ قیمت برای یک مگابایت از دیسک ۵۰ برابر کمتر از DRAM بود.

بنابراین دیسک‌های مغناطیسی و حافظه اصلی سه فرق عمده با هم دارند: دیسک‌ها غیر فرارند، به دلیل اینکه مغناطیسی هستند. سرعت دیسک‌ها کمتر است چون قطعات مکانیکی دارند و قیمت دیسک‌ها کمتر است چون هزینه تولید آنها در مقایسه با DRAM ها پایین‌تر است.

۴-۱- مدارهای مجتمع

تکنولوژی ساخت پردازنده‌ها و حافظه‌ها در طول سالیان همواره در حال پیشرفت بوده است و به همین دلیل پردازنده‌ها و حافظه‌ها همواره در حال بهبود بوده‌اند. جدول ۱-۱ تکنولوژی‌هایی را نشان می‌دهد که در سالهای مختلف مورد استفاده قرار گرفته‌اند. در این جدول یک تخمین از افزایش سرعت به ازای یک هزینه واحد آورده شده است.

جدول ۱-۱: افزایش سرعت در هزینه واحد برای تکنولوژیهای مختلف

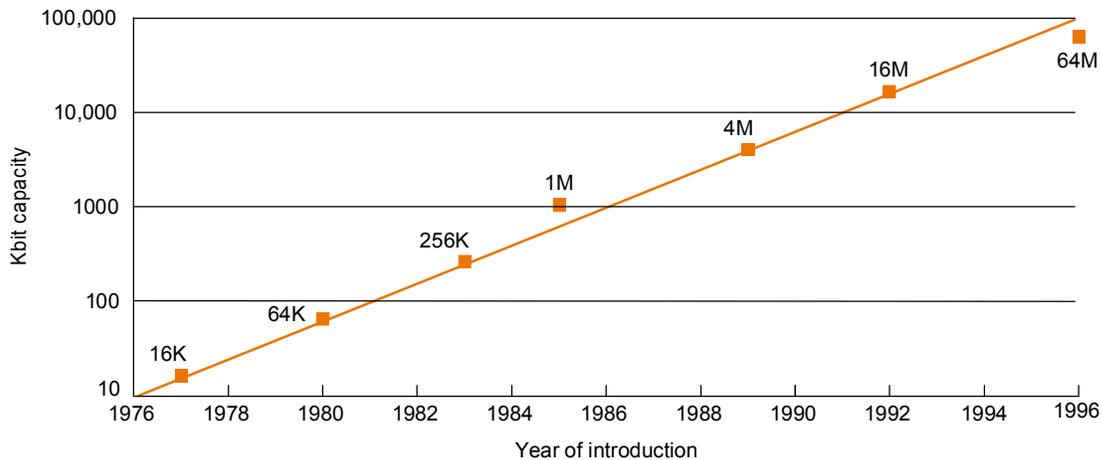
Year	Technology used in computers	Related Performance / unit cost
1951	Vacuum tub	1
1965	Transistor	35
1975	Integrated Circuit	900
1995	Very large-scale Integrated circuit	2,400,000

این بخش در مورد تکنولوژی که سرعت کامپیوتر را از سال ۱۹۷۵ تحت تأثیر قرار داده و در سالهای آتی نیز این اتفاق خواهد افتاد، صحبت خواهد نمود. چون تکنولوژی نشان می‌دهد که کامپیوترها چه کارهایی را می‌توانند انجام دهند و با چه سرعتی در حال پیشرفت می‌باشند، بنابراین ما اعتقاد داریم که همه متخصصین کامپیوتر باید با اصول مدارهای مجتمع آشنا شوند.

یک ترانزیستور به زبان ساده یک سوئیچی است که دارای دو وضعیت روشن و خاموش (on و off) بوده و با جریان برق کنترل می‌شود. یک مدار مجتمع حدوداً ده‌ها تا صدها عدد از این ترانزیستورها را در داخل یک چیپ ترکیب می‌نماید برای نشان دادن نرخ افزایش سرسام‌آور تعداد ترانزیستورها از صدها عدد به میلیون‌ها عدد در داخل یک چیپ به آنها VLSI^۱ گفته می‌شود.

شکل ۱-۵ میزان افزایش ظرفیت DRAM ها را از سال ۱۹۷۷ به بعد نشان می‌دهد. همان طور که مشاهده می‌شود ظرفیت DRAM ها تقریباً هر سه سال، چهار برابر شده است که در نتیجه در طول بیست سال ظرفیت DRAM ها حدوداً ۱۶۰۰۰ برابر شده است! این میزان افزایش قابل توجه در سرعت و ظرفیت مدارهای مجتمع، طراحی سخت‌افزار و نرم‌افزار را تحت تأثیر قرار می‌دهد و همین عاملی است که فهمیدن مدارهای مجتمع را ضروری می‌سازد.

^۱ - Very Large Scale Integrated



شکل ۱-۵: میزان افزایش ظرفیت DRAM ها در طول زمان. واحد محور y ها Kbits می باشد.

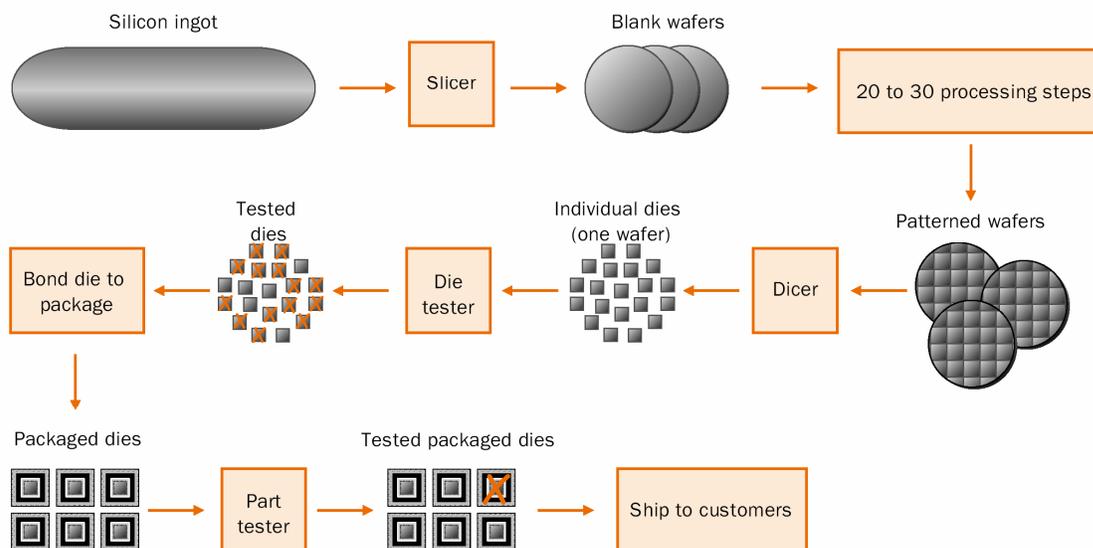
حال که ارتباط افزایش حجم حافظه ها را با تکنولوژی متوجه شدیم در این قسمت می خواهیم روند ساخته شدن یک چیپ را مورد بررسی قرار دهیم. ساختن یک چیپ از سیلیکون شروع می شود (سیلیکون ماده ای است که در سنگ یافت می شود و جزو عناصر گروه چهارم جدول تناوبی مندلیف می باشد). چون سیلیکون جریان الکتریسیته را به خوبی عبور نمی دهد به آن عنصر نیمه هادی گفته می شود. با استفاده از یک پروسه شیمیایی مخصوص این امکان وجود دارد که موادی به سیلیکون اضافه نماییم تا در سطح آن محل هایی ایجاد کنیم که در آن محل ها بتوانیم یکی از سه قطعه زیر را بسازیم:

- مواد هادی جریان الکتریسیته (شبهه به مس و آلومینیوم)
- مواد عایق (مانند پلاستیک و شیشه)
- فضایی (area) که تحت شرایط خاصی می تواند هادی و یا عایق (مانند یک کلید) باشد

ترانزیستورها در طبقه آخر قرار می گیرند. یک مدار VLSI شامل میلیونها هادی، عایق و سوئیچ می باشد که بر روی یک فضای کوچک سیلیکونی ساخته می شوند. پروسه ساخت مدارهای مجتمع بر روی هزینه چیپ ها تأثیر زیادی دارد و بنابراین برای طراحان کامپیوتر اهمیت زیادی دارد. شکل ۱-۶ این پروسه ساخت را نشان می دهد. این پروسه از یک شمش سیلیکونی شروع می شود. سپس یک شمش برش خورده و به قطعات کوچکتر و نازکتری به نام Wafer تقسیم می شود. این ویفرها سپس چند مرحله پروسه را طی می کنند، که در آن پروسه ها یک سری مواری شیمیایی بر روی هر ویفر قرار می گیرد که ترانزیستورها، هادیها و عایقها را ایجاد می کنند، بعد از طی کردن این مراحل ویفرها برش داده می شوند و به قطعات کوچکتری به نام die تبدیل می شوند. بعد از این مرحله die های خراب کنار گذاشته می شوند. تعداد die های سالم به کل die های یک ویفر yield نامیده می شود. بعد از تست die ها، هر کدام از die های سالم در داخل یک package، به پین های ورودی و خروجی وصل می شود این

پروسه bonding نامیده می‌شود. Die‌های package شده دوباره تست می‌شوند چون ممکن است هنگام package کردن خرابی‌هایی اتفاق بیفتد. پس از تست شدن، package به مشتری فروخته می‌شود. توجه: یکی از مسائل جدید در طراحی کامپیوتر توان تلفاتی می‌باشد. اهمیت توان اتلافی فقط برای کاربردهای قابل حمل (Portable) که طول عمر باتری در آنها مهم است (مانند کامپیوترهای کیفی) نیست بلکه برای کامپیوترهای رومیزی که نیز با افزایش کلاک اهمیت مسأله توان بیشتر می‌شود. پردازنده Alpha 21264 در فرکانس 600 مگاهرتز به طور شگفت‌انگیزی فقط 72 وات توان مصرف می‌کند. توان اتلافی یکی از مسائلی است که می‌تواند سرعت پردازنده‌ها را محدود کند.

Building Computer Chips



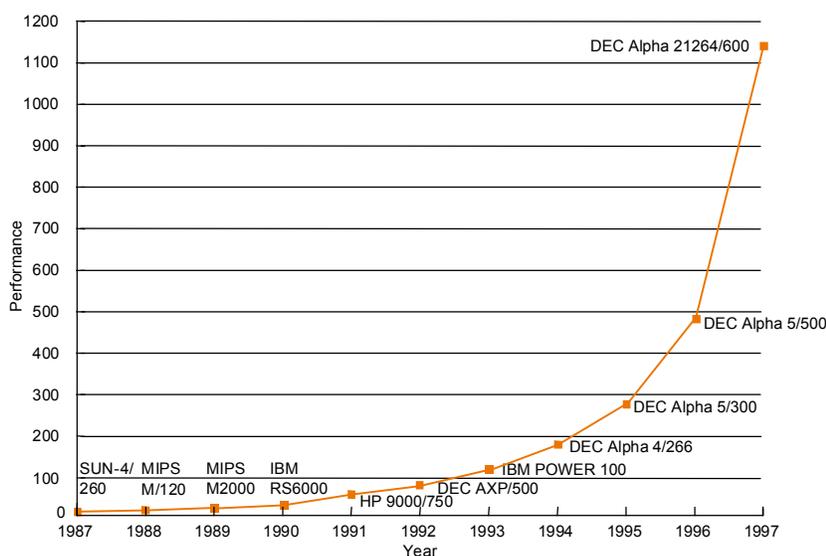
شکل ۱-۶: پروسه ساخت چیپ

۱-۵- باورهای غلط (fallacies) و اشتباهات (Pitfalls)

هدف از داشتن یک بخش تحت این عنوان که در آخر هر فصل وجود دارد این است که مطالبی را توضیح دهد که ممکن است شما آنها را اشتباه متوجه شده باشید: ما این باور غلط را fallacy نامگذاری خواهیم کرد. وقتی که در مورد یک باور غلط صحبت می‌کنیم برای توضیح دادن آن یک مثال خواهیم آورد. همچنین ما در مورد اشتباهات (Pitfalls) صحبت خواهیم کرد. اغلب این اشتباهات

تعمیم دادن اصولی است که فقط در محدوده کوچکتري درست می‌باشند. هدف از این بخش این است که به شما کمک کند تا جلوی اشتباهات خود را در ماشین‌هایی که استفاده و یا طراحی می‌کنید بگیرید. **اشتباه:** چشم پوشی کردن از پیشرفت سخت‌افزار هنگامی که یک ماشین جدید را طراحی می‌کنیم اشتباه است.

فرض کنید که شما می‌خواهید یک ماشین را در مدت ۳ سال تولید کنید و ادعا دارید که این ماشین یک فروشنده فوق‌العاده خواهد بود چون سرعت آن ۳ برابر سرعت ماشین‌های دیگر که امروزه وجود دارند خواهد بود. متأسفانه ماشین شما احتمال دارد فروشنده ضعیفی باشد به دلیل اینکه افزایش سرعت برای صنعت (افزایش سرعت سخت‌افزار) ممکن است به ماشین‌هایی منجر شود که سرعتشان مثل سرعت ماشین شما باشد. به طور مثال فرض کنید ۵۰٪ رشد سالیانه سرعت داشته باشیم (به دلیل پیشرفت تکنولوژی ساخت سخت‌افزار)، در این صورت برای ماشینی که امروزه سرعت آن x می‌باشد می‌توان انتظار داشت که پس از ۳ سال سرعتی در حدود $1.5^3 x = 3.4x$ خواهد داشت. بنابراین ماشین شما مزیت سرعتی نخواهد داشت! بسیاری از پروژه‌ها در شرکت‌های کامپیوتری به این دلیل تعطیل شدند که یا این قانون را در نظر نگرفتند (پیشرفت تکنولوژی ساخت سخت‌افزار) و یا اینکه دیرتر از موعد خودشان تمام شدند و سرعت ماشینی که با تأخیر به بازار آمد از میانگین سرعت صنعت پایین‌تر بود. این واقعیت ممکن است در همه صنایع وجود داشته باشد ولی بهبود سریع هزینه بر سرعت (Cost / Performance) در صنعت کامپیوتر آن را به یک امر بسیار مهم تبدیل کرده است. شکل ۱-۷ میزان افزایش سرعت پردازنده‌ها در طول زمان را نشان می‌دهد.



شکل ۱-۷: افزایش سرعت کامپیوترها بین سالهای ۱۹۸۷ تا ۱۹۹۷

۱-۶- موضوع واقعی: ساختن چیپ‌های پنتیوم (Pentium)

در انتهای هر فصل این کتاب، بخشی تحت عنوان موضوع واقعی (Real stuff) وجود دارد که مطالب موجود در همان فصل را به کامپیوترهایی که در عمل از آنها استفاده می‌کنید، مرتبط می‌سازد. این بخش‌ها همیشه تکنولوژی‌هایی را که برای کامپیوترهای IBM PC مورد استفاده قرار می‌گیرد توضیح خواهند داد و اغلب تکنولوژی Apple Macintosh را نیز بحث خواهند کرد. برای این فصل ما مفاهیم مدارهای مجتمع را به چیپ‌هایی که در IBM PC ها استفاده می‌شوند، مرتبط خواهیم نمود.

۱-۷- صحبت پایانی

هر چند پیشگویی این امر مشکل است که واقعاً بهبود هزینه بر سرعت کامپیوترها در آینده چگونه خواهد بود، ولی آنها مطمئناً بهتر از کامپیوترهای فعلی خواهند بود. برای مشارکت در این پیشرفت، طراحان کامپیوتر و برنامه‌نویسان باید مسائل زیادی را بفهمند و یاد بگیرند.

هر جفت طراحان سخت‌افزار و نرم‌افزار سیستم‌های کامپیوتری را به صورت لایه‌ای طراحی می‌کنند که در آن لایه‌های پایین‌تر جزئیات را از دید لایه‌های بالاتر مخفی نگه می‌دارند. این اصول سطح تجرد برای فهمیدن سیستم‌های کامپیوتری، امروزه ضروری است. شاید مهمترین مثال از سطح تجرد، واسط بین سخت‌افزار و نرم‌افزار سطح پایین می‌باشد که ISA نامیده می‌شود. با فرض ثابت نگه داشتن ISA پیاده‌سازی‌های مختلفی از یک معماری با هزینه و سرعت متفاوت را ممکن می‌سازد. همه این پیاده‌سازی‌ها نرم‌افزار یکسانی را اجرا می‌کنند. از طرف دیگر ممکن است این معماری مانع پیشرفت گردد و نیاز به این داشته باشد که واسط مورد نظر تغییر کند. تکنولوژی‌های کلیدی برای پردازنده‌های مدرن، سیلیکون و کامپایلر می‌باشند. به طور واضح شما باید ویژگیهای سیلیکون و کامپایلر را درک نمائید. چیزی که اهمیت آن کمتر از فهمیدن مفهوم مدارهای مجتمع نیست فهمیدن سرعت تغییرات مورد انتظار تکنولوژی‌های ساخت IC می‌باشد. یک مثال مرتبط با این مسأله این است که سرعت DRAM ها هر سه سال ۴ برابر می‌شود. در حالی که سیلیکون عامل مؤثری در پیشرفت سریع سخت‌افزار می‌باشد، ایده‌های جدید در طراحی کامپیوترها، هزینه بر سرعت کامپیوترها را بهبود داده است. دو مورد از ایده‌های کلیدی، به کارگیری موازی سازی در کامپیوترها و سلسله مراتب حافظه می‌باشد. نوعاً موازی سازی از طریق پایپلاین و به کارگیری مفهوم دسترسی محلی به حافظه از طریق حافظه‌های کش می‌باشد.

راهنمای فصول این کتاب

همان طور که توضیح داده شد پنج قسمت کلاسیکی برای کامپیوتر وجود دارد: datapath، بخش کنترل، حافظه، ورودی و خروجی این پنج قسمت به عنوان یک قالب کاری (frame work) برای قسمت‌های باقی مانده این کتاب مورد استفاده قرار می‌گیرد:

- مسیر داده (datapath): فصول ۴ و ۵ و ۶
- بخش کنترل: فصول ۵ و ۶
- حافظه: فصل ۷
- ورودی: فصل ۸
- خروجی: فصل ۸

فصل ۶ توضیح خواهد داد که چگونه پایپلاین در پردازنده موازی‌سازی انجام می‌دهد و فصل ۷ در مورد اینکه چگونه سلسله مراتب حافظه خاصیت محلی بودن را بکار می‌گیرد، بحث خواهد نمود. بقیه فصل‌ها مقدمات و مباحث تکمیلی در مورد این موضوعات را مطرح می‌کنند. فصل ۲ در مورد سرعت پردازنده‌ها صحبت خواهد کرد و بنابراین توضیح می‌دهد که چگونه یک کامپیوتر را ارزیابی کنیم. فصل ۳ در مورد مجموعه دستورات که همان واسط بین کامپایلرها و ماشین است صحبت خواهد کرد و تأکید خواهد نمود که کامپایلرها و زبانهای برنامه‌نویسی چه نقشی در استفاده از خواص مجموعه دستورات خواهند داشت.

فصل ۹ این کتاب با بحثی در مورد سیستم‌های چند پردازنده‌ای این کتاب را به پایان خواهد برد.

۸-۱- چشم‌انداز تاریخی و مطالعه اضافی

در این کتاب در آخر هر فصل بخشی تحت این عنوان وجود دارد. در این بخش‌ها ما ممکن است توسعه یک ایده را از طریق یک‌سری از ماشین‌ها بررسی کنیم یا بعضی از پروژه‌های مهم را توضیح دهیم. ما در هر مورد مراجع را معرفی خواهیم نمود و اگر علاقمند بودید می‌توانید به آنها مراجعه کنید و بیشتر آنها را ارزیابی کنید. در حالت کلی مطالب این بخش در مورد مطالب همان فصل مربوطه خواهد بود و مباحث تاریخی و تکمیلی را در مورد مطالب آن فصل بیان خواهد نمود.

۹-۱- کلمات و اصطلاحات کلیدی

لیستی از کلمات و اصطلاحات در آخر هر فصل و در پیوست آورده شده است. این کلمات ایده‌های کلیدی بحث شده در آن فصل و پیوست را منعکس می‌کند. اگر شما معنای کلمه و اصطلاح گفته شده

در این قسمت را هنوز نفهمیده‌اید، می‌توانید به واژگان کلمات در آخر کتاب مراجعه کنید. همه کلمات و اصطلاحات کلیدی در این قسمت توضیح داده شده است.

ارزیابی کارآیی پردازنده‌ها

این فصل توضیح خواهد داد که چگونه کارایی^۶ (سرعت) یک کامپیوتر را اندازه بگیریم و آن را گزارش کنیم. در این فصل سعی خواهیم نمود که به سؤالات زیر جواب دهیم:

- چرا کارایی مهم است؟
- چگونه ما می‌توانیم کارایی را به دقت تعریف کنیم؟
- چگونه طراحی سخت‌افزار کارایی نرم‌افزار را تحت تأثیر قرار می‌دهد؟
- چگونه در دنیای واقعی کارایی اندازه گرفته شود؟
- چرا بعضی از سخت‌افزارها بهتر از بقیه عمل می‌کنند؟
- کدام یک از فاکتورهای کارایی به سخت‌افزار وابسته‌اند؟
- ما برای اجرا شدن سریع‌تر یک برنامه به یک ماشین جدید نیاز داریم یا نیازمند یک سیستم-عامل جدید هستیم؟
- مجموعه دستورات یک ماشین چگونه کارایی را تحت تأثیر قرار می‌دهند؟

در این فصل عاملهای تعیین کننده در کارایی کامپیوترها مورد بحث قرار خواهد گرفت. بررسی کردن کارایی به این دلیل مهم است که کارایی سخت‌افزار اغلب کلیدی‌ترین اثر را در کل سیستم متشکل از سخت‌افزار و نرم‌افزار برعهده دارد. ما معمولاً در این فصل کارایی و سرعت را به جای هم استفاده می‌کنیم. ولی در بعضی از موارد این دو عبارت کاملاً معادل هم نیستند. به عنوان مثال اگر هدف ما استفاده از ماشینی باشد که تعداد کار بیشتری را در واحد زمان انجام دهد، ماشینی که سریعتر است بعضی مواقع نمی‌تواند هدف ما را برآورده کند و ما مجبور هستیم در این مواقع از ماشینی استفاده کنیم که قابلیت‌های بیشتری داشته باشد (کارتر باشد). ما در این کتاب اغلب از اصطلاح کارایی استفاده خواهیم کرد.

اظهار نظر کردن در مورد کارایی یک سیستم واقعاً کار مشکلی است. سیستم‌های بزرگ نرم‌افزاری با جزئیات زیاد به همراه بازه وسیعی از تکنیک‌های افزایش سرعت که توسط طراحان سخت‌افزار به کار گرفته می‌شوند، از جمله مواردی هستند که صحبت کردن درباره کارایی را مشکل می‌نمایند.

تقریباً کار غیر ممکن است که برگه راهنمایی از مجموعه دستورات یک کامپیوتر را بردارید و با یک برنامه مشخص، تعیین نمایید که آن کامپیوتر برنامه شما را چقدر سریعتر اجرا خواهد کرد. در واقع

⁶ - Performance

برای برنامه‌های کاربردی مختلف، ممکن است مترهای مختلفی مناسب باشد و قسمت‌های مختلف یک کامپیوتر ممکن است نقش متفاوتی در تعیین سرعت کل سیستم داشته باشند.

البته در انتخاب بین کامپیوترهای مختلف، کارآیی تقریباً همیشه یک مشخصه مهم به حساب می‌آید. اندازه‌گیری و مقایسه ماشین‌های مختلف برای خریداران و در نتیجه برای طراحان یک امر حیاتی محسوب می‌شود که این مطلب را فروشندگان کامپیوترها به خوبی درک می‌کنند. اغلب، فروشندگان دوست دارند که شما کامپیوترهای آنها را خیلی ساده نگاه کنید و این نگاه کردن ممکن است بدون توجه به نیاز خریدار، که برنامه‌های کاربردی خودش را بر روی آن اجرا خواهد کرد انجام بگیرد. در بعضی مواقع ادعاهایی در مورد کامپیوترها می‌شود که دید مفیدی برای بعضی از برنامه‌های کاربردی در اختیار قرار نمی‌دهند. بنابراین فهمیدن چگونگی اندازه‌گیری کارآیی یک کامپیوتر و محدودیت‌های اندازه‌گیری در انتخاب یک ماشین، بسیار مهم می‌باشند.

علاقمندی ما نسبت به کارآیی کامپیوتر ماورای این است که فقط کارآیی کامپیوتر را از بیرون آن اندازه بگیریم. ما نیاز به این داریم که بفهمیم چه چیزهایی کارآیی یک کامپیوتر را تحت تأثیر قرار می‌دهند. فهمیدن اینکه چرا بعضی قسمت‌های نرم‌افزار همان طور که ما انتظار داریم کار نمی‌کنند، چرا یک مجموعه از دستورات می‌تواند به گونه‌ای پیاده‌سازی شود که بهتر عمل کند، و اینکه چگونه بعضی از قابلیت‌های سخت‌افزار کارآیی را تحت تأثیر قرار می‌دهند، همگی از مسائلی است که ما علاقمند به پیدا کردن جواب مناسب برای آنها هستیم.

به طور مثال برای بهبود سرعت یک سیستم نرم‌افزاری، ممکن است احتیاج داشته باشیم که فاکتورهای جدید سخت‌افزاری را که تأثیر بر کل سیستم دارند را بشناسیم و میزان اهمیت این فاکتورها را درک نمائیم. این فاکتورها ممکن است شامل موارد زیر باشد: برنامه‌های کاربران کدام دستورات ماشین را بیشتر استفاده می‌کنند، سخت‌افزار چگونه این مجموعه دستورات را پیاده‌سازی می‌کند و سیستم‌های IO و حافظه چقدر خوب کار می‌کنند؟ فهمیدن اینکه چگونه تأثیر این فاکتورها را بر کارآیی تعیین کنیم، برای درک انگیزه‌های پشت صحنه طراحی بعضی از قابلیت‌های ویژه ماشین که ما آنها را در فصل‌های آتی خواهیم دید، خیلی مهم می‌باشند. این فصل روش‌های اندازه‌گیری کارآیی را توضیح خواهد داد. در بخش ۲-۲ ما معیارهای اندازه‌گیری کارآیی را از دو دیدگاه استفاده‌کننده کامپیوتر و طراح سخت‌افزار بیان خواهیم کرد. در بخش ۲-۳ مشاهده خواهیم کرد که چگونه این معیارها به هم وابستگی دارند و فرمولی جهت اندازه‌گیری کارآیی ارائه خواهد شد که در سراسر این کتاب از آن استفاده خواهد شد. بخش‌های ۲-۴ و ۲-۵ در مورد نحوه انتخاب benchmark هایی که برای ارزیابی

کارایی کامپیوترها استفاده می شود صحبت خواهد کرد. همچنین یاد می گیریم که چگونه سرعت اجرای دسته‌ای از برنامه‌ها را به طور دقیق خلاصه و جمع‌بندی نمائیم و گزارش تهیه کنیم. بخش ۲-۶ یک دسته از benchmark های معمول استفاده شده برای CPUها را توضیح خواهد داد و با استفاده از این benchmark ها سرعت چند پردازنده اینتل را اندازه خواهد گرفت. سرانجام در فصل ۲-۷ ما به دسته-ای از کج فهمی‌ها و دام‌هایی که معمولاً طراحان و تحلیل کننده‌های کارایی را گرفتار می‌کنند، خواهیم پرداخت.

تعریف کارایی

وقتی که گفته می‌شود کارایی این کامپیوتر بالاتر از آن یکی است، منظور چیست؟ هر چند این سؤال ساده به نظر می‌رسد ولی یک مقایسه با هواپیماهای مسافربری نشان می‌دهد که این سؤال چقدر قابل تأمل است.

شکل ۲-۱ تعدادی از هواپیماهای مسافربری را به همراه سرعت متوسط آنها، برد مسافتی و ظرفیتشان نشان می‌دهد. اگر بخواهیم بینیم که کدام یک از هواپیماهای موجود در جدول کارایی بالاتری دارد، اول باید کارایی را تعریف کنیم. به طور مثال، با در نظر گرفتن اندازه‌گیری‌های مختلف، می‌بینیم که هواپیمایی که بالاترین سرعت متوسط را دارد Concorde، و هواپیمایی که بیشترین ظرفیت را دارد 747 می‌باشد.

Plane	DC to Paris	Speed	Passengers	Throughput (pmp)
747	6.5 hours	610 mph	470	286,700
Concorde	3 hours	1350 mph	132	178,200

شکل ۲-۱: مقایسه کارایی هواپیماهای مختلف

بنابراین اگر ما بخواهیم کارایی این دو هواپیما را مقایسه کنیم. به دو صورت می‌توانیم این کار را انجام دهیم. می‌توانیم به این صورت تعریف کنیم که هواپیمای کاراتر هواپیمایی است که یک مسافر را از نقطه‌ای به نقطه دیگر در کمترین زمان جابجا می‌کند. پس اگر شما علاقمند باشید که زودتر به مقصد برسید واضح است که Concorde کاراتر می‌باشد. روش دوم مقایسه، برای این مثال تعداد مسافر جابجا شده در واحد زمان است. بر اساس این معیار همان طور که در ستون آخر جدول نیز نشان داده شده است 747 کاراتر می‌باشد.

به طور مشابه ما می‌توانیم کارایی کامپیوترها را به طرق مختلف تعریف نمائیم. اگر شما یک برنامه را بر روی دو کامپیوتر مختلف اجرا کنید، خواهید گفت که کامپیوتری سریعتر است که آن برنامه را سریعتر اجرا می‌کند. ولی اگر برنامه خود را بر روی یک ترمینال در مرکز کامپیوتری اجرا کنید که دارای دو کامپیوتر است که به طور مشترک توسط چندین نفر استفاده می‌شوند در این صورت خواهید گفت که کامپیوتری سریعتر است که تعداد زیادتری برنامه را در روز اجرا کند. همان طور که یک کاربر دوست دارد که زمان پاسخ (response time) یا همان زمان اجرا (execution time) که مدت زمان بین شروع و خاتمه یک کار است، بر روی یک کامپیوتر برای برنامه او کمتر شود، در یک مرکز کامپیوتر، مدیر مرکز اغلب دوست دارد که تعداد کار انجام گرفته در واحد زمان (throughput) افزایش پیدا کند.

مثال

برای نشان دادن نحوه استفاده از ایده‌های جدید، مثالهای ویژه‌ای در سراسر این کتاب آورده شده است. ما ابتدا مثال را می‌آوریم و بعد به آن پاسخ می‌دهیم. سعی کنید که مثال را خودتان جواب دهید، اگر نتوانستید و یا به پاسخ خود اطمینان نداشتید جواب را مشاهده نمائید. مثالهایی که آورده می‌شود شبیه به مسائلی است که در آخر هر فصل آورده می‌شود. این اولین مثال از این کتاب است:

مثال: آیا تغییرات انجام شده زیرین بر روی یک سیستم کامپیوتری، throughput را زیاد می‌کند؟ زمان پاسخ دهی را کاهش می‌دهد؟ و یا هر دو؟

۱. جایگزین کردن پردازنده سیستم با یک نسخه سریعتر

۲. اضافه کردن یک پردازنده جدید به سیستمی که از چندین پردازنده برای کارهای مختلف استفاده می‌کند. به طور مثال سیستمی که عملیات رزواسیون یک خطوط هواپیمایی را انجام می‌دهد.

جواب: کم کردن زمان پاسخ دهی تقریباً همیشه throughput را بهبود می‌دهد. بنابراین در مورد اول، زمان پاسخ دهی و throughput هر دو بهتر می‌شوند. در مورد ۲، برنامه هیچ کسی سریعتر انجام نمی‌گیرد، بنابراین فقط throughput بهتر می‌شود. اگر در مورد ۲، نیاز به زمان پاسخ دهی سریع همانند افزایش throughput اهمیت داشته باشد، سیستم ممکن است مجبور کند که درخواستها در داخل یک صف قرار بگیرند. در این حالت افزایش throughput، زمان پاسخ دهی را نیز کاهش خواهد داد چون زمان انتظار را در داخل صف کاهش خواهد داد. بنابراین در تعداد زیادی از سیستم‌های واقعی کامپیوتری، تغییر زمان پاسخ دهی یا throughput، اغلب دیگری را نیز تحت تأثیر قرار می‌دهد.

ما در مبحث کارایی کامپیوترها، در ابتدا در طول چند فصل بر روی زمان پاسخ دهی تمرکز خواهیم کرد. در فصل ۸ که در ارتباط با سیستم‌های ورودی/خروجی می‌باشد ما اندازه‌گیری‌های مربوط به throughput را نیز صحبت خواهیم کرد.

برای اضافه کردن سرعت (کارایی)، ما می‌خواهیم که زمان پاسخ دهی یا زمان اجرا را کاهش دهیم، بنابراین می‌توان زمان اجرا و کارایی را برای یک ماشین X به صورت زیر مرتبط کرد:

۱

$$\text{کارایی ماشین X} = \text{-----}$$

زمان اجرا بر روی ماشین X

فرمول فوق به صورت زیر نیز نوشته می‌شود:

$$performance_x = \frac{1}{(ExecutionTime)_x}$$

این به این معنی است که برای دو ماشین X و Y، اگر کارایی X بیشتر از کارایی Y باشد داریم:

$$\begin{aligned} performance_x &> performance_y \\ \Rightarrow \frac{1}{(ExecutionTime)_x} &> \frac{1}{(ExecutionTime)_y} \\ \Rightarrow (ExecutionTime)_y &> (ExecutionTime)_x \end{aligned}$$

یعنی اینکه اگر X سریعتر از Y باشد زمان اجرا بر روی Y بیشتر از زمان اجرا بر روی X خواهد بود. در مبحث طراحی کامپیوتر، ما اغلب سرعت دو کامپیوتر مختلف را به صورت عددی مقایسه می‌کنیم. مثلاً می‌گوئیم ماشین X، n بار سریعتر از ماشین Y می‌باشد و می‌نویسیم.

$$\frac{performance_x}{performance_y} = n$$

اگر X، n بار سریعتر از Y باشد در این صورت زمان اجرا بر روی Y، n بار بیشتر از زمان اجرا بر روی X خواهد شد:

$$\frac{performance_x}{performance_y} = \frac{(ExecutionTime)_y}{(ExecutionTime)_x} = n$$

مثال: اگر ماشین A یک برنامه را در ۱۰ ثانیه و ماشین B آن را در ۱۵ ثانیه اجرا کند، ماشین A چند برابر از B سریعتر است؟

جواب: طبق فرمول بالایی می‌دانیم که اگر ماشین A، n بار سریعتر از B باشد می‌توان نوشت:

$$\frac{performance_A}{performance_B} = \frac{(ExecutionTime)_B}{(ExecutionTime)_A} = n$$

بنابراین نسبت کارایی به صورت زیر حساب می‌شود:

$$n = \frac{15}{10} = 1.5$$

بنابراین A، 1.5 برابر سریعتر از B می‌باشد.

در مثال بالایی می‌توان گفت که ماشین B، ۱/۵ برابر کندتر از ماشین A عمل می‌کند. چون داریم:

$$\frac{performance_A}{performance_B} = 1.5 \Rightarrow \frac{performance_A}{1.5} = performance_B$$

برای سادگی در مقایسه کارایی کامپیوترها ما از عبارت سریعتر بودن استفاده خواهیم کرد. به دلیل اینکه کارایی و زمان اجرا معکوس هم می‌باشند، افزایش کارایی نیازمند این است که زمان اجرا را کاهش دهیم. برای جلوگیری از گیج شدن بین جمله‌های افزایش و کاهش، ما معمولاً وقتی که منظورمان "افزایش کارایی" و "کاهش زمان اجراست"، می‌گوئیم "بهبود کارایی" یا "بهبود زمان اجرا".^۷

۲-۲- اندازه‌گیری کارایی

اکنون زمان اندازه‌گیری کارایی است، کامپیوتری که همان مقدار کار را در زمان کمتری انجام می‌دهد سریعتر است. زمان اجرای برنامه با ثانیه اندازه گرفته می‌شود. اما زمان می‌تواند بسته به اینکه ما چه چیزی را شمارش می‌کنیم، به گونه‌ای دیگر تعریف گردد. سراسرترین تعریف زمان، زمان پاسخ‌دهی یا زمان انقضا می‌باشد.^۹ این اصطلاحات به معنی کل زمان لازم برای خاتمه یک کار است، که این زمان شامل دسترسی به حافظه دیسک، دسترسی به حافظه اصلی، دسترسی به IOها، سربار اجرای سیستم-عامل و هر زمان دیگری است.

⁷ - Performance improvement

⁸ - Execution time improvement

⁹ - Elapsed time

کامپیوترها اغلب به صورت اشتراک زمانی استفاده می‌شوند و یک پردازنده ممکن است بر روی چندین برنامه به صورت همزمان کار کند. در این صورت سیستم ممکن است تلاشش در این راستا باشد که به جای کاهش زمان اجرای یک برنامه، throughput را کاهش دهد. بنابراین ما می‌خواهیم که بین زمان انتظار و زمانی که cpu فقط برای ما کار می‌کند تفاوت قائل شویم. زمان اجرای cpu (cpu execution time) یا به طور ساده cpu time مدت زمانی است که cpu صرف کار ما می‌کند و شامل زمان انتظار IO یا زمان اجرای برنامه‌های دیگر نمی‌شود. (به خاطر بیاورید که زمان پاسخ‌دهی به کاربر زمان انقضای برنامه است نه زمان cpu time).

cpu time می‌تواند به دو بخش تقسیم گردد: زمانی از cpu که صرف اجرای خود برنامه می‌شود (user cpu time) و زمانی از cpu که صرف اجرای سیستم‌عامل در ارتباط با این برنامه می‌شود. (معمولاً system cpu time).

سیستم عامل unix دستوری به نام time دارد که می‌تواند زمان انقضا را نشان دهد به طور مثال این دستور می‌تواند نتیجه‌ای مانند عبارت زیر را برگرداند:

90.7u 12.95 2:39 65%

در عبارت فوق زمان user cpu time 90.7 ثانیه، زمان system cpu time 12.9 ثانیه، زمان انقضا 2 دقیقه و 39 ثانیه (۱۵۹ ثانیه) و درصدی از زمان انقضا که مربوط به زمان cpu time می‌باشد به صورت زیر است:

$$\frac{90.7+12.9}{159} = 0.65$$

یا 65 درصد. در این مثال بیش از $\frac{1}{3}$ از زمان انقضا صرف انتظار برای IO، اجرای برنامه‌های دیگر یا هر دو شده است.

بعضی مواقع وقتی صحبت از زمان اجرای cpu می‌شود ما معمولاً از زمان system cpu time چشم‌پوشی می‌کنیم و این معمولاً به دلیل غیر دقیق بودن اندازه‌گیری سیستم‌عامل که خودش را اندازه می‌گیرد و نیز به دلیل بی‌عدالتی در مواقعی است که می‌خواهیم system cpu time دو ماشین که سیستم‌عاملهای مختلفی را اجرا می‌کنند مقایسه کنیم. از طرف دیگر، کدهای سیستمی در بعضی ماشینها، جزو کدهای کاربری ماشین دیگر می‌باشد. تقریباً هیچ برنامه‌ای نمی‌تواند بدون حضور سیستم‌عامل بر روی سخت‌افزار اجرا شود، بنابراین زمان اجرای برنامه را می‌توان مجموع user cpu time و system cpu time در نظر گرفت.

ما بین کارآیی محاسبه شده بر اساس زمان انقضا و براساس زمان cpu execution time فرق قائل خواهیم شد. ما عبارت کارآیی سیستم (system performance) را برای زمان انقضا برای یک سیستم خالی از برنامه به کار می‌بریم و عبارت cpu performance را در ارتباط با user cpu time به کار خواهیم برد. ما در این فصل بر روی کارآیی cpu (cpu performance) تمرکز خواهیم کرد، هر چند بحث ما برای سرعت می‌تواند برای زمان انقضا و برای cpu time نیز مورد استفاده قرار گیرد. هر چند که ما مانند کاربران کامپیوتر به زمان دقت می‌کنیم، وقتی ما جزئیات یک ماشین را بررسی می‌کنیم راحت‌تر این است که معیارهای دیگری را برای سرعت در نظر بگیریم. علی‌الخصوص طراحان کامپیوتر ممکن است در مورد یک ماشین به این صورت فکر کنند که با استفاده از یک معیاری نشان دهند که سخت‌افزار آن کامپیوتر عملیات پایه‌ای را چقدر سریعتر انجام می‌دهد. تقریباً همه کامپیوترها در ساختارشان از یک کلاک (clock) استفاده می‌کنند که این کلاک سیگنالی است که به صورت پیرودیک تولید می‌شود و مشخص می‌کند اتفاقات در سخت‌افزار کی رخ دهند. به هر پیرودی کلاک یک سیکل کلاک (clock cycle, ticks, clock ticks, clock periods, clocks, cycles) گفته می‌شود. طراحان معمولاً طول پیرودی کلاک را به صورت زمان یک سیکل کامل کلاک (مانند ۲ نانو ثانیه) و یا به صورت نرخ کلاک (مانند ۵۰۰ مگاهرتز یا 500 mhz) که معکوس پیرودی کلاک است به کار می‌برند. در بخش بعدی ما ارتباط بین سیکل کلاک و زمان اجرای برنامه کاربر را توضیح خواهیم داد.

۲-۳- ارتباط بین معیارها (مترها)

کاربران و طراحان برای اندازه‌گیری سرعت از معیارهای مختلفی استفاده می‌کنند. اگر ما بتوانیم این معیارها را به هم مرتبط سازیم در این صورت می‌توانیم اثر یک تغییر در طراحی بر روی سرعت کامپیوتر را که توسط کاربر مشاهده می‌شود، مشخص نماییم. فرمول زیر نحوه محاسبه زمان اجرای cpu را نشان می‌دهد. این فرمول ارتباط کلاک با زمان cpu را نشان می‌دهد.

$$\text{پیرودی کلاک} \times \text{تعداد کلاک‌های لازم برای اجرای برنامه} = \text{زمان اجرای cpu}$$

چون فرکانس کلاک (clock rate) و پیرودی کلاک معکوس هم هستند فرمول بالایی را می‌توان به صورت زیر نیز نمایش داد:

$$\text{تعداد کلاک‌های لازم برای اجرای برنامه}$$

$$= \text{زمان اجرای cpu} \text{ -----}$$

فرکانس کلاک

این فرمول به وضوح نشان می‌دهد که طراح سخت‌افزار می‌تواند با کاهش طول پریود کلاک یا تعداد کلاک‌های لازم برای اجرای یک برنامه، سرعت اجرای برنامه را افزایش دهد (سرعت را بهبود دهد). ما در ادامه این فصل و در فصول ۵ و ۶ و ۷ خواهیم دید که طراح نمی‌تواند همزمان هم طول پریود کلاک را کاهش دهد و هم تعداد کلاک‌های لازم برای برنامه را کاهش دهد و اغلب باید یک بالانس بین آنها برقرار کند (trade off). اکثر تکنیک‌هایی که تعداد کلاک‌ها را کاهش می‌دهند معمولاً طول پریود کلاک را افزایش می‌دهند.

مثال: برنامه دلخواه ما بر روی کامپیوتر A که فرکانس کلاک آن ۴۰۰ مگاهرتز است در ۱۰ ثانیه اجرا می‌شود. ما می‌خواهیم به یک طراح کامپیوتر کمک کنیم که یک ماشین به نام B بسازد که برنامه ما را در ۶ ثانیه انجام دهد. طراح کامپیوتر به ما گفته است که افزایش فرکانس کلاک امکان‌پذیر است ولی این افزایش فرکانس طراحی بخش‌های دیگر cpu را تحت تأثیر قرار خواهد داد و باعث خواهد شد که ماشین B برای اجرای این برنامه تعداد کلاک‌هایی در حدود ۱/۲ برابر ماشین A نیاز داشته باشد. ما چه فرکانس کلاکی را برای رسیدن به این هدف از طراح بخواهیم؟

جواب: در ابتدا تعداد کلاک‌های لازم برای اجرای برنامه بر روی ماشین A را حساب می‌کنیم:

تعداد کلاک‌های لازم برای اجرای برنامه بر روی A

$$\text{فرکانس کلاک ماشین A} = \frac{\text{زمان اجرای برنامه بر روی A}}{\text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی A}}$$

فرکانس کلاک ماشین A

تعداد کلاک‌های لازم برای اجرای برنامه بر روی A

$$10 \text{ s} = \frac{\text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی A}}{400 \times 10^6}$$

$$400 \times 10^6$$

$$\text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی A} = 4000 \times 10^6$$

زمان اجرای برنامه توسط ماشین B به صورت زیر است:

تعداد کلاک‌های لازم برای اجرای برنامه بر روی B

$$\text{فرکانس کلاک ماشین B} = \frac{\text{زمان اجرای برنامه بر روی B}}{\text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی B}}$$

فرکانس کلاک ماشین B

$$1.2 \times (\text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی A})$$

$$1.2 \times 4000 \times 10^6$$

$$6 \text{ s} = \frac{1.2 \times (\text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی A})}{1.2 \times 4000 \times 10^6}$$

$$\text{فرکانس کلاک ماشین B} = 800 \text{ MHz}$$

بنابراین ماشین B برای اینکه برنامه ما را در ۶ ثانیه اجرا کند باید فرکانسش دو برابر فرکانس ماشین A باشد.

تلاقی سخت‌افزار و نرم‌افزار (Hardware Software Interface)

در سراسر این کتاب شما بخشهایی تحت این عنوان خواهید دید. این بخش‌ها تأثیر متقابل بعضی ویژگی‌های نرم‌افزار (نرم‌افزار می‌تواند یک برنامه، یک کامپایلر و یا یک سیستم‌عامل باشد) و بعضی ویژگی‌های سخت‌افزار بر همدیگر را بیان خواهد نمود و آن را به صورت کاملاً واضح و روشنی توضیح خواهد داد. این بخشها باعث خواهند شد که ما به خاطر بسپاریم که طراحی سخت‌افزار و نرم‌افزار در بسیاری از موارد تأثیر متقابلی بر هم دارند.

در معادلات موجود در مثالهای قبلی جایی برای تعداد دستورات مورد نیاز یک برنامه در نظر گرفته نشده بود. اما به هر حال چون کامپایلر برای اجرا شدن برنامه دستورهایی به زبان ماشین تولید می‌کند و ماشین نیز برای اجرای برنامه باید آن دستورات را اجرا نماید، بنابراین زمان اجرای یک برنامه باید به تعداد دستورات تولید شده برای یک برنامه وابسته باشد. یکی از راههایی که می‌توان برای محاسبه زمان اجرا در نظر گرفت این است که بگوئیم زمان اجرا مساوی است با تعداد دستورات یک برنامه ضربدر زمان متوسط لازم برای اجرای یک دستور. بنابراین تعداد سیکل‌های کلاک مورد نیاز برای یک برنامه می‌تواند به صورت زیر محاسبه شود:

(تعداد کلاک متوسط لازم برای اجرای هر دستور) × (تعداد دستورات برنامه) = تعداد کلاک‌های لازم برای اجرای یک برنامه
عبارت تعداد کلاک‌های متوسط مورد نیاز برای هر دستور^{۱۰} اغلب CPI نامیده می‌شود. چون دستورات مختلف بسته به اینکه چه کاری انجام می‌دهند، ممکن است زمان اجرای متفاوتی داشته باشند، CPI متوسط زمان اجرای همه دستوراتی است (زمان اجرا به کلاک) که در داخل یک برنامه قرار دارند. CPI روشی برای مقایسه دو پیاده‌سازی مختلف از یک مجموعه دستورات را نیز در اختیار قرار می‌دهد چون در این دو پیاده‌سازی تعداد دستورات مورد نیاز برنامه مطمئناً یکی خواهد بود.

¹⁰ - clock cycles per instruction

مثال: فرض کنید ما دو پیاده‌سازی از یک مجموعه دستورات واحد در اختیار داشته باشیم. ماشین A برای یک برنامه دارای پیوند کلاک ۱ نانوثانیه بوده و CPI آن ۲ می‌باشد و ماشین B دارای پیوند کلاک ۲ نانوثانیه بوده و CPI آن ۱/۲ می‌باشد. برای این برنامه کدام ماشین سریعتر است و چقدر؟
جواب: می‌دانیم تعداد دستوراتی که هر کدام از ماشین‌ها برای این برنامه اجرا می‌کنند یکی است. ما این تعداد دستورات را I می‌نامیم. در ابتدا تعداد کلاکهای پردازنده را برای هر کدام از ماشین‌ها حساب می‌کنیم:

$$I \times 2.0 = \text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی ماشین A}$$

$$I \times 1.2 = \text{تعداد کلاک‌های لازم برای اجرای برنامه بر روی ماشین B}$$

حال زمان اجرا برنامه را برای هر کدام از ماشین‌ها حساب می‌کنیم:

$$\begin{aligned} \text{پیوند کلاک} \times \text{تعداد کلاک‌های لازم برای اجرای برنامه} &= \text{زمان اجرای برنامه بر روی ماشین A} \\ &= I \times 2.0 \times 1 \text{ ns} = 2 \times I \text{ ns} \end{aligned}$$

همین طور برای ماشین B داریم:

$$I \times 1.2 \times 2 \text{ ns} = 2.4 \times I \text{ ns} = \text{زمان اجرای برنامه بر روی ماشین B}$$

واضح است که ماشین A سریعتر است چون زمان اجرای پایین‌تری دارد.. برای اینکه بدانیم A چقدر سریعتر از B است به صورت زیر عمل می‌کنیم:

$$\frac{\text{کارایی ماشین A}}{\text{کارایی ماشین B}} = \frac{\text{زمان اجرای برنامه بر روی ماشین B}}{\text{زمان اجرای برنامه بر روی ماشین A}} = \frac{2.4 \times I \text{ ns}}{2 \times I \text{ ns}} = 1.2$$

بنابراین ماشین A برای این برنامه، 1.2 برابر سریعتر از ماشین B می‌باشد.

مثال: آیا موارد زیر می‌توانند به تنهایی نشان دهنده کارایی باشند؟

۱. تعداد کلاکهای لازم برای یک برنامه

۲. تعداد دستورات یک برنامه

۳. پیوند کلاک

۴. تعداد کلاک لازم برای اجرای یک دستور

۵. تعداد متوسط دستوراتی که در یک ثانیه اجرا می‌شوند.

جواب: هیچکدام از موارد نمی‌توانند به تنهایی نشان دهنده کارایی باشند. کارایی با زمان اجرا ارتباط معکوس دارد و باید زمان اجرا به نوعی در معیار مشخص شده باشد. در مورد اول تعداد کلاک معیار

مناسبی نیست چون مشخص نیست که پریود کلاک چقدر است و بنابراین نمی‌توانیم زمان اجرا را بدست آوریم. در مورد دوم تعداد دستورات نیز معیار مناسبی نیست چون نوع و زمان اجرای دستورات مشخص نیست و نمی‌توانیم زمان اجرای برنامه را بدست آوریم. دستورات مختلف زمان اجرای متفاوتی دارند به عنوان مثال یک دستور ضرب نسبت به یک دستور جمع زمان اجرای بالاتری دارد و دستورات ممیز شناور که با داده‌های ممیزدار کار می‌کنند نسبت به دستوراتی که از داده‌های نوع صحیح استفاده می‌کنند، زمان اجرای بالاتری دارند. در مورد سوم پریود کلاک نیز معیار مناسبی نیست چون مشخص نیست که برنامه به چند کلاک برای اجرا نیاز دارد و نمی‌توانیم زمان اجرای برنامه را مشخص کنیم. در مورد چهارم تعداد کلاک لازم برای اجرای هر دستور نیز معیار مناسبی نیست چون تعداد دستورات و پریود کلاک مشخص نشده و نمی‌توانیم زمان اجرا را حساب کنیم. در مورد پنجم تعداد متوسط دستوراتی که در یک ثانیه اجرا می‌شود نیز به تنهایی معیار مناسبی نیست چون نوع دستورات و تعداد کلاک لازم برای هر دستور و فرکانس کلاک مشخص نیست و نمی‌توانیم زمان اجرا را بدست آوریم.

مثال: اگر دو ماشین ISA یکسان داشته باشد کدام یک از موارد زیر در دو ماشین همیشه برابر می‌باشد؟

۱. فرکانس کلاک

۲. CPI

۳. زمان اجرا

۴. تعداد دستورات

۵. MIPS

جواب: اگر دو ماشین مختلف وجود داشته باشند که ISA یکسانی را پیاده سازی کرده باشند یعنی مجموعه دستوراتی که پشتیبانی می‌کنند مثل هم باشد، در این صورت حتماً تعداد دستورات یک برنامه که برای این دو ماشین کامپایل بشود در هر دو کامپایل یکی خواهد بود.

مثال: یک طراح کامپایلر می‌خواهد بین دو قطعه کد برای یک ماشین مشخص یکی را انتخاب نماید. بر اساس پیاده سازی سخت افزار، سه کلاس مختلف از دستورات وجود دارد: کلاس A، کلاس B و کلاس C که به ترتیب برای اجرا شدن نیاز به ۱، ۲ و ۳ کلاک دارند. قطعه کد اول دارای ۵ دستور است: ۲ مورد از کلاس A، ۱ مورد از کلاس B و ۲ مورد از کلاس C و قطعه کد دوم دارای ۶ دستور است: ۴ مورد از کلاس A، ۱ مورد از کلاس B و ۱ مورد از کلاس C.

الف) کدام قطعه کد سریعتر اجرا می‌شود و چقدر؟

ب) CPI را برای هر قطعه کد حساب کنید؟

جواب: الف) برای بدست آوردن تعداد کلاک از فرمول زیر استفاده می‌کنیم:

$$\text{تعداد کلاکهای لازم} = \sum_{i=1}^n (CPI_i \times C_i)$$

که در آن CPI_i تعداد کلاکهای لازم برای اجرای دستورات از نوع کلاس i می‌باشد و C_i تعداد دستورات کلاس i می‌باشد.

$$\text{کلاک } 1 = 2 \times 1 + 1 \times 2 + 2 \times 3 = 10$$

$$\text{کلاک } 2 = 4 \times 1 + 1 \times 2 + 1 \times 3 = 9$$

چون هر دو قطعه کد بر روی یک ماشین اجرا می‌شوند و در نتیجه پریود کلاک برای هر دو یکی است بنابراین قطعه کدی سریعتر است که تعداد کلاک کمتری لازم داشته باشد. بنابراین در این مثال با اینکه قطعه کد ۲ تعداد دستورات بیشتری دارد ولی سریعتر اجرا می‌شود.

ب) برای بدست آوردن CPI از فرمول زیر استفاده می‌کنیم:

تعداد کلاکها

$$CPI = \frac{\text{تعداد کلاکها}}{\text{تعداد دستورات}}$$

تعداد دستورات

۱۰

$$CPI = \frac{10}{5} = 2$$

۵

۹

$$CPI = \frac{9}{6} = 1.5$$

۶

مثال: فرض کنید یک ماشین داشته باشیم که با سرعت ۵۰۰ MHz کار کند و فرض کنید که این ماشین دارای سه کلاس دستور از نوع A و B و C باشد که به ترتیب برای اجرا شدن به یک، دو و سه کلاک نیاز داشته باشند. می‌خواهیم دو کامپایلر را بر روی این دو ماشین امتحان کنیم. این کامپایلرها برای تولید کد از یک برنامه بزرگ مورد استفاده قرار می‌گیرند.

کامپایلر اول به تعداد ۵ میلیارد دستور از نوع A، یک میلیارد دستور از نوع B و یک میلیارد دستور از نوع C استفاده می‌کند. و کامپایلر دوم به تعداد ۱۰ میلیارد دستور از نوع A، یک میلیارد دستور از نوع B و یک میلیارد دستور از نوع C استفاده می‌کند.

الف) کدام قطعه کد از لحاظ زمان اجرا سریعتر است؟

ب) کدام قطعه کد از لحاظ تعداد MIPS سریعتر است؟

جواب: MIPS به عنوان یک معیار برای مقایسه کارایی پردازنده‌ها مورد استفاده قرار می‌گیرد. MIPS که مخفف عبارت Millions of Instruction Per second می‌باشد برای یک کامپیوتر به این صورت تعریف می‌شود: تعداد میلیون دستوری که یک کامپیوتر در یک ثانیه انجام می‌دهد. این مثال نشان می‌دهد که MIPS معیار مناسبی برای کارایی نیست چون در بعضی مواقع نمی‌تواند زمان اجرا را به درستی نشان بدهد.

الف) برای زمان اجرا از فرمول زیر استفاده می‌کنیم:

$$\begin{aligned} \text{زمان اجرا برای قطعه کد اول} &= \left(\sum_{i=1}^n CPI_i \times C_i \right) \times \text{پریود کلاک} \\ &= (5 \times 1 + 1 \times 2 + 1 \times 3) \times 10^9 \times 10 \times 10^9 \\ &= 10 \times 10^6 \times \text{پریود کلاک} = \frac{10 \times 10^6}{500 \times 10^6} = 20 \text{ s} \end{aligned}$$

$$\begin{aligned} \text{زمان اجرای قطعه کد دوم} &= \left(\sum_{i=1}^n CPI_i \times C_i \right) \times \text{پریود کلاک} \\ &= (10 \times 1 + 1 \times 2 + 1 \times 3) \times 10^9 = 15 \times 10^9 \times 10 \times 10^9 \\ &= \frac{15 \times 10^9}{500 \times 10^6} = 30 \text{ s} \end{aligned}$$

پس از نظر زمان اجرا قطعه کد اول سریعتر است.

ب) برای MIPS از فرمول زیر استفاده می‌کنیم:

تعداد دستورات

$$MIPS = \frac{\text{تعداد دستورات}}{\text{زمان اجرا}}$$

$10^6 \times$ زمان اجرا

$$(5 + 1 + 1) \times 10^9$$

$$\text{MIPS} = \frac{\text{قطعه کد اول}}{20 \times 10^6} = 350$$

$$(10 + 1 + 1) \times 10^9$$

$$\text{MIPS} = \frac{\text{قطعه کد دوم}}{30 \times 10^6} = 400$$

بنابراین از نظر MIPS قطعه کد دوم سریعتر است.

مثال: تعداد دستورات برنامه زیر مشخص کنید.

```
Li $a0, 1000 // a0 = 1000
```

```
Loop: sub $a0, $a0, 1 // a0 = a0 - 1
```

```
bne $a0, $0, Loop // if a0 != 0 go to loop
```

جواب: تعداد دستورات یک برنامه را می‌توان به دو صورت شمارش نمود: استاتیکی و دینامیکی. تعداد دستورات استاتیکی در حقیقت همان تعداد خطوط یا تعداد دستورات ظاهری برنامه است ولی تعداد دستورات دینامیکی تعداد دستوراتی است که در عمل اجرا می‌شوند. بنابراین برای این مثال داریم:

$$3 = \text{تعداد دستورات استاتیک}$$

چون تعداد دستورات ظاهری برنامه ۳ می‌باشد.

$$2001 = \text{تعداد دستورات دینامیک}$$

توضیحات هر دستور در کنار همان دستور بعد از علامت // آورده شده است. اگر این برنامه را trace کنیم می‌بینیم که حلقه برنامه به تعداد ۱۰۰۰ مرتبه اجرا می‌شود. چون تعداد دستورات داخل حلقه ۲ است پس تعداد کل دستوراتی که در این حلقه اجرا خواهد شد ۲۰۰۰ مورد می‌شود. توجه داریم که یک دستور نیز در خارج حلقه قرار دارد. پس در کل تعداد دستورات دینامیکی برنامه که در عمل اجرا می‌شوند، ۲۰۰۱ مورد خواهد شد.

مثال: دو پروسور زیر را که دارای ISA یکسان یا سازگار (ISA compatible) می‌باشند از نظر کارایی باهم مقایسه کنید؟

۱. یک پردازنده ۸۰۰ مگاهرتز AMD Duron که CPI آن ۱/۲ است

۲. یک پردازنده ۱ گیگاهرتز پنتیوم III (P3) که CPI آن ۱/۵ است

جواب: چون دو پردازنده ISA یکسانی دارند پس برای یک برنامه فرضی تعداد دستورات برای هر جفت آنها یکی است و داریم:

$$\begin{aligned} \text{پریود کلاک} \times \text{CPI} \times \text{تعداد دستورات} &= \text{زمان اجرای برنامه روی AMD} \\ &= I \times 1/2 \times 1/25 \text{ ns} = 1/5 \text{ ns} \times I \end{aligned}$$

$$\text{P3 روی برنامه اجرای زمان} = I \times 1/5 \times 1 \text{ ns} = 1/5 \text{ ns} \times I$$

چون زمان اجرای برنامه فرضی بر روی دو پردازنده یکی است پس کارایی آنها یکسان است! مثال: دو کامپیوتر زیر را که کلاک مختلف ولی ISA یکسان دارند با هم مقایسه کنید.

۱. P4 2.5GHz

۲. P4 3GHz

جواب: چون دو کامپیوتر ISA یکسانی دارند پس تعداد دستورات یک برنامه فرضی بر روی آنها برابر است و همچنین چون دو کامپیوتر از یک نسل می‌باشند (هر دو از نسل P4 هستند) از نظر ساختاری مانند هم بوده و CPI آنها نیز یکی است و داریم:

$$\text{پریود کلاک} \times \text{CPI} \times \text{تعداد دستورات} = \text{زمان اجرا}$$

$$0.4 \text{ ns} \times \text{CPI} \times \text{تعداد دستورات} = \text{زمان اجرای ۱}$$

$$0.33 \text{ ns} \times \text{CPI} \times \text{تعداد دستورات} = \text{زمان اجرای ۲}$$

$$\begin{array}{r} 0.4 \\ \text{زمان اجرای ۱} \\ \text{کارایی ۲} \end{array} = \frac{\text{کارایی ۱}}{\text{زمان اجرای ۲}} = \frac{0.33}{0.4} = 1/2$$

پس پردازنده P4 3GHz به اندازه ۱/۲ برابر سریعتر از پردازنده P4 2.5GHz است.

مثال: مقایسه دو کامپایلر برای یک ماشین

فرض کنید دو کامپایلر برای یک ماشین داشته باشیم: کامپایلر معمولی و کامپایلر بهینه. کامپایلری که بهتر عمل می‌کند کدهایی تولید می‌کند که تعداد دستورات آن ۵٪ (پنج درصد) کمتر است و همچنین CPI آنها ۱۰٪ پایین‌تر است. این دو کامپایلر را مقایسه کنید.

جواب: فرض کنید کامپایلر معمولی را کامپایلر پایه در نظر گرفته و تعداد دستورات تولید شده توسط آن را برای یک برنامه فرضی با I ، CPI کد تولید شده توسط آن را با CPI_B و پریود کلاک ماشین را با T نشان دهیم در این صورت داریم:

پریود کلاک $\times CPI_B \times$ تعداد دستورات پایه = زمان اجرای کد تولید شده توسط کامپایلر پایه

$$= I \times CPI_B \times T$$

$CPI \times T$ = تعداد دستورات = زمان اجرای کد تولید شده با کامپایلر بهینه

$$= (0.95 \times I) \times (0.9 \times CPI_B) \times T$$

$$= 0.86 \times (I \times CPI_B \times T)$$

$$= 0.86 \times \text{زمان اجرای کامپایلر پایه}$$

$$\frac{1/0.86}{1/0.9} = \frac{\text{زمان اجرای پایه}}{\text{کارایی کامپایلر بهینه}}$$

$$= 1.17$$

$$\frac{0.86}{0.9} = \frac{\text{کارایی کامپایلر پایه}}{\text{زمان اجرای بهینه}}$$

توجه: با توجه به مثالهای متنوعی که حل کردیم می‌توان نتیجه گرفت که روشهای متعددی برای افزایش سرعت پردازنده‌ها وجود دارد. بعضی از این روشها مربوط به تکنیک‌های کامپایلری و نرم-افزاری است که می‌توان توسط این تکنیکها کدهای با تعداد دستورات کم و با CPI پایین تولید نمود (CPI می‌تواند بر اساس نوع و تعداد دستورات از کدی به کد دیگر تغییر کند در واقع ترکیب دستورات^{۱۱} CPI را عوض می‌کند). بعضی دیگر از این روشها سخت‌افزاری هستند و می‌توان سخت-افزارهایی طراحی نمود که پریود کلاک پایینی داشته باشند و تعداد کلاک لازم برای اجرای هر دستور نیز پایین باشد (در واقع با طراحی سخت‌افزارهای مناسب می‌توان CPI را کم کرد).

قانون امدال

قانون امدال در معماری کامپیوتر قانونی است که می‌توان به کمک آن ارزش کارهایی که قرار است انجام شوند را از قبل پیش بینی نمود. این کارها عمدتاً مربوط به بهبودهایی است که در سخت‌افزار پردازنده‌ها داده می‌شود. این قانون به شرح زیر است:

¹¹ - Instruction mix

قانون امدال: مقدار بهبودهایی که انجام می‌شود به میزان مؤثر بودن آنها محدود می‌شود.

برای توضیح این قانون سیستم واحدی دانشگاه را در نظر می‌گیریم. در این سیستم بعضی دروس یک واحدی، بعضی دو واحدی، بعضی سه واحدی و بعضی نیز چهار واحدی‌اند. همان‌طور که می‌دانیم در محاسبه معدل ترم، دروسی بیشتر تأثیر دارند که تعداد واحد آنها بالاتر است. بنابراین در این سیستم ارزش دروس چهار واحدی بیشتر از سه واحدی، ارزش دروس سه واحدی بیشتر از دو واحدی، ارزش دروس دو واحدی نیز بیشتر از یک واحدی است. بنابراین ما بهتر است بیشترین وقت و سرمایه‌گذاری را بر روی دروسی انجام دهیم که بیشترین تأثیر را در معدل دارند. در معماری کامپیوتر قانون امدال می‌گوید: قبل از اینکه بهبودی را انجام دهی یا کاری را انجام دهی باید به میزان مؤثر بودنش نیز توجه کنی و آن کار را وقتی انجام دهی که تأثیرش در مقابل کاری که انجام می‌دهی قابل توجه باشد. در قالب فرمول قانون امدال به صورت زیر بیان می‌شود:

زمان اجرایی که تحت تأثیر قرار می‌گیرد

زمان اجرایی که تحت تأثیر قرار نمی‌گیرد + ----- = زمان اجرای برنامه بعد از بهبود

میزان بهبود

مثال: فرض کنید که تصمیم گرفته باشیم سخت افزار پردازنده را تغییر دهیم تا سرعت اجرای دستورات floating point دو برابر شود. دو برابر کردن سرعت یک ایده خوب است ولی باید ببینیم در قبال کاری که انجام می‌دهیم چقدر بهبود در کل زمان اجرای یک برنامه به وجود می‌آید. به عبارتی باید ببینیم دستورات floating point چه میزان در برنامه ما نقش دارند. فرض کنید برای یک برنامه فرضی فقط ۱۰ درصد زمان اجرای برنامه (T) شامل دستورات floating point باشد. برای این برنامه تأثیر بهبود سخت‌افزار را بررسی کنید.

جواب: با استفاده از قانون امدال می‌توان نوشت:

$$0.10 T$$

$$0.95 T = 0.90 T + \text{-----} = \text{زمان اجرای برنامه بعد از بهبود}$$

۲

یعنی زمان اجرای کل برنامه فقط ۵ درصد بهبود خواهد داشت. این میزان بهبود خیلی بالا نیست بنابراین بر طبق قانون امدال دو برابر کردن سرعت اجرای دستورات floating point به امید بهبود زمان اجرای برنامه در این مثال ارزش فنی ندارد.

نتیجه قانون امدال: قسمتهایی از برنامه را بهبود دهیم که بیشتر استفاده می‌شوند چون بیشترین تأثیر را در کل زمان اجرای برنامه دارند. قانون امدال بیان می‌کند که موارد پر استفاده را بهبود دهیم و یا به عبارتی می‌گوید: *Make the common case fast*.

۲-۴- انتخاب برنامه‌ها برای اندازه‌گیری کارایی

کاربر کامپیوتری که هر روز مجموعه برنامه‌های ثابتی را اجرا می‌کند، یک کاندیدای خوب برای ارزیابی سرعت یک کامپیوتر جدید است. مجموعه برنامه‌هایی که توسط این کاربر اجرا می‌شوند، یک workload نامیده می‌شوند. برای مقایسه کارایی دو سیستم کامپیوتری این کاربر به سادگی زمان اجرای یک workload را بر روی دو ماشین مقایسه می‌کند. اما مشکلی که در اینجا وجود دارد این است که اکثر کاربران در این شرایط قرار ندارند و باید روشهای دیگری را برای ارزیابی کارایی استفاده کنند و امیدوار باشند که این روشها بتوانند نشان دهنده کارایی workload های واقعی باشند. این روشها معمولاً مبتنی بر استفاده از مجموعه‌ای از benchmark ها می‌باشند. Benchmark ها مجموعه‌ای از برنامه‌ها هستند که برای اندازه‌گیری کارایی انتخاب می‌شوند. Benchmark ها یک workload را تشکیل می‌دهند که کاربر امیدوار است که نشان دهنده کارایی یک workload واقعی باشد.

یک benchmark خوب باید بتواند کارایی برنامه‌های واقعی را منعکس نماید و علی‌الخصوص benchmark باید دارای ترکیب واقعی از دستورات باشد (یعنی شامل ترکیبی از دستورات باشد که در برنامه‌های واقعی ظاهر می‌شوند). امروزه ثابت شده است که بهترین benchmark ها برنامه‌های واقعی هستند. تعدادی از benchmark های معمول عبارتند از:

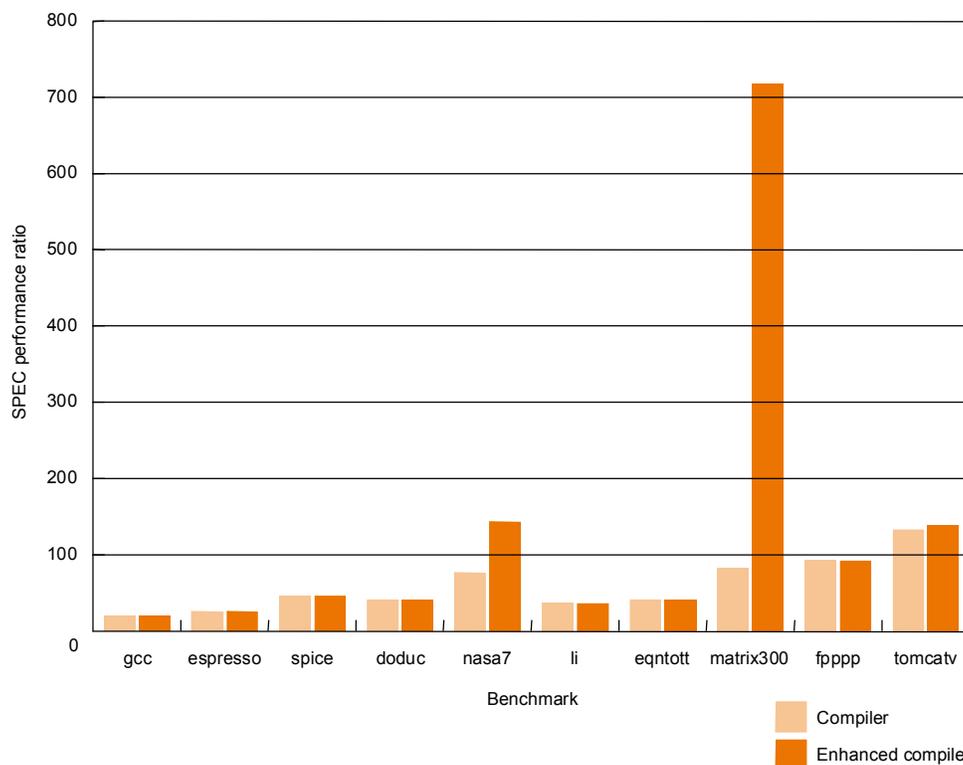
- Adobe Photoshop for image processing
- BABCO SYSmark for office applications
- Unreal Tournament 2003 for 3D games
- SPEC (System Performance Evaluation Cooperative)

امروزه برخی از سازنده‌های کامپیوتر بر اساس اطلاعاتی که از کاربران و benchmark های آنان دارند، روشهای نادرستی را در پیش می‌گیرند و اطلاعات غلطی را به کاربران می‌دهند. این سازندگان یا از تکنیکهای کامپایلری استفاده می‌کنند که برای یک benchmark مشخص کدهای بهینه‌ای تولید می‌کنند که بر روی سخت افزار سریعتر اجرا می‌شوند و یا اینکه از تکنیکهای سخت‌افزاری استفاده می‌کنند که یک ترکیب مشخص از دستورات را که مربوط به یک benchmark مشخص است را با سرعت بالاتری اجرا می‌کنند. استفاده از برنامه‌های واقعی به عنوان benchmark استفاده از روشهای فوق برای افزایش

سرعت اجرا شدن benchmark ها را تقریباً غیر ممکن می‌سازد حتی اگر بهبودی نیز حاصل شود این بهبود چون مربوط به برنامه‌های واقعی است، قابل استفاده خواهد بود.

استفاده کردن از benchmark هایی که اندازه آنها کوچک است و یا کارایی آنها به قسمت کوچکی از کد برنامه بستگی دارد خطر بهینه سازی‌های کامپایلری یا سخت افزاری را به همراه خواهد داشت. به طور مثال مجموعه benchmark های SPEC در ابتدا برای این منظور انتخاب شدند که از برنامه‌های واقعی برای اندازه‌گیری کارایی استفاده کنند. متأسفانه اولین نسخه از مجموعه SPEC که در سال ۱۹۸۹ وارد بازار شد، دارای یک benchmark به نام matrix300 بود که فقط دارای یک سری عملیات ضرب ماتریسی بود. در واقع ۹۹ درصد از زمان اجرا مربوط به یک خط از این برنامه بود. این واقعیت که زمان خیلی زیادی از اجرای برنامه مربوط به یک خط از برنامه بود که در آن یک خط، یک عملیات مشخص انجام می‌گرفت، باعث شد تا چند کمپانی کامپیوتری کامپایلرهای را تولید و یا خریداری کنند تا زمان اجرای این benchmark را بهینه سازی کنند. شکل ۲-۳ نسبت کارایی^{۱۲} (معکوس زمان اجرا) را برای یک ماشین با دو کامپایلر مختلف نشان می‌دهد. برنامه‌های این شکل مجموعه benchmark های SPEC89 می‌باشند. در این شکل برای هر برنامه دو میله وجود دارد که میله سمت راستی مربوط به کامپایلر بهینه‌سازی شده برای matrix300 می‌باشد. همان طور که مشاهده می‌شود کامپایلر بهینه-سازی شده برای matrix300 اثری بر روی ۸ مورد از ۱۰ مورد از برنامه‌های مجموعه benchmark های SPEC89 ندارد، اما کارایی matrix300 را بیش از ۹ برابر بهبود داده است. کاربری که براساس matrix300 قضاوت می‌کند حسابی دچار اشتباه بزرگی خواهد شد به دلیل اینکه بهبود داده شده فقط مربوط به matrix300 بوده و برای برنامه‌های دیگر کاربردی نخواهد داشت.

¹² - Performance ratio



شکل ۲-۳: نسبت کارایی مجموعه benchmark های SPEC89 برای یک ماشین با دو کامپایلر مختلف ضعف مربوط به matrix300 باعث شد که این benchmark در نسخه بعدی مجموعه SPEC یعنی SPEC92 کنار گذاشته شود.

حال سؤال این است که چرا همه مردم از برنامه‌های واقعی برای اندازه‌گیری کارایی استفاده نمی‌کنند؟ یکی از دلایل این است که benchmark های کوچک به علت کوچک بودن به راحتی کامپایل و شبیه‌سازی می‌شوند (حتی بعضی مواقع دستی کامپایل می‌شوند) و به همین دلیل در شروع طراحی جذابیت‌های خاصی دارند. مخصوصاً چون در شروع طراحی یک ماشین هنوز کامپایلری برای آن نوشته نشده است، این benchmark ها جایگاه خاصی دارند. دلیل دیگری که برای استفاده از benchmark های کوچک وجود دارد این است که آنها راحت‌تر از برنامه‌های بزرگ استاندارد سازی می‌شوند و بنابراین نتایج منتشر شده زیادی برای benchmark های کوچک وجود.

اگرچه استفاده از benchmark های کوچک در مراحل اولیه طراحی قابل توجیه است، اما هیچ دلیل قانع کننده‌ای وجود ندارد که از آنها برای ارزیابی کارایی سیستم‌های کامپیوتری که وارد بازار شده‌اند استفاده کنیم. در گذشته یکی از دلایلی که از benchmark های کوچک به جای برنامه‌های بزرگ استفاده می‌شد این بود که آنها بر روی ماشین‌های مختلف به راحتی پورت می‌شدند و از طریق آنها می‌شد سیستم‌های مختلف را مقایسه نمود. اما امروزه این دلیل خیلی درست نیست و استفاده از

benchmark های کوچک فقط در مراحل اولیه طراحی مجاز است و بعد از آن باید از برنامه‌های واقعی برای ارزیابی کارایی ماشین‌ها استفاده شود.

پس از آنکه ما مجموعه مناسبی از benchmark ها را انتخاب نموده و توسط آنها کارایی را اندازه‌گیری کردیم، می‌توانیم یک گزارش برای کارایی تهیه کنیم. توضیحاتی که در گزارش آورده می‌شود باید طوری باشد که اندازه‌گیری ما قابل تکرار باشد. گزارش ما باید شامل یک لیست از همه چیزهایی باشد که ممکن است یک فرد دیگر برای تکرار آزمایش‌های ما به آنها نیاز داشته باشد. این لیست باید شامل نسخه سیستم عامل، کامپایلرها، ورودی‌ها و همچنین نحوه پیکربندی ماشین باشد. به طور مثال توصیف سیستمی^{۱۳} ماشینی که برای بدست آوردن نتایج شکل ۲-۳ از آن استفاده شده است در شکل ۲-۴ نشان داده شده است

Hardware	
Model number	Powerstation 550
CPU	41.67 MHz POWER 4164
FPU	Integrated
Number of CPUs	1
Cache size per CPU	64K data / 8K instruction
Memory	64 MB
Disk subsystem	2 400-MB SCSI
Network interface	NA
Software	
OS type and rev	AIX v3.1.5
Compiler rev	AIX XL C/6000 Ver. 1.1.5 AIX XL Fortran Ver. 2.2
Other software	None
File system type	AIX
Firmware level	NA
System	
Tuning parameters	None
Background load	None
System state	Multiuser (single-user-login)

شکل ۲-۴: توصیف سیستمی ماشین استفاده شده برای اندازه‌گیری‌های شکل ۲-۳

¹³ - System description

Synthetic های Benchmark

یک synthetic benchmark، یک برنامه است که فقط برای منظور اندازه‌گیری کارایی نوشته شده است و معمولاً اندازه آن کوچک بوده و به راحتی بر روی CPU های مختلف پورت می‌شود. بنابراین برای مقایسه سیستم‌ها راحت‌تر می‌باشند. هدف از طراحی synthetic benchmark ها این بوده که یک برنامه واحد داشته باشیم که ویژگیهای تعداد زیادی از برنامه‌ها را داشته باشد و تعداد تکرار دستورات در داخل برنامه متناسب با تعداد تکرار دستورات در تعداد زیادی از benchmark ها باشد. معروفترین benchmark های synthetic عبارتند از Whetstone و Dhrystone .

چون این نوع benchmark ها واقعی نیستند نشان دهنده هیچ چیز خاصی نیستند. به طور مثال هیچ خروجی در اختیار کاربر قرار نمی‌دهند که کاربر بداند برنامه به درستی اجرا شده است و نیز اگر بر روی یک ماشین کارایی آنها بهبود پیدا کرد هیچ دلیلی وجود ندارد که برای برنامه‌های واقعی نیز چنین باشد و کارایی آن ماشین برای سایر برنامه‌ها نیز بهبود یابد. اندازه این benchmark ها کوچکتر است و بنابراین به راحتی تکنیکهای بهینه سازی کامپایلر و سخت‌افزار برای آنها قابل اعمال است. نکته: بنا به دلایل فوق synthetic benchmark ها برای اندازه‌گیری کارایی انتخاب مناسبی نیستند.

۲-۵- مقایسه کردن و خلاصه سازی کارایی

حال که یاد گرفتیم برنامه‌هایی را به عنوان benchmark انتخاب نمائیم و برای مقایسه از معیارهای زمان اجرا و throughput استفاده کنیم، ممکن است تصور کنید که مقایسه کردن کارایی کار سر راستی باشد. اما معمولاً به جای یک benchmark از چند benchmark استفاده می‌شود و ما باید یاد بگیریم که چگونه کارایی یک گروه از benchmark ها را اندازه‌گیری کنیم. اگر چه خلاصه کردن یک مجموعه از اندازه‌گیری‌ها اطلاعات کمی در اختیار قرار می‌دهد، ولی بازار و حتی کاربران اغلب ترجیح می‌دهند که یک عدد ساده از کارایی داشته باشند و با استفاده از آن مقایسه‌های خود را انجام دهند. حال سؤال کلیدی این است که این عدد چگونه محاسبه می‌شود؟ شکل ۲-۵ زمان اجرای دو برنامه را بر روی دو ماشین مختلف نشان می‌دهد.

	Computer A	Computer B
Program 1 (seconds)	1	10
Program 2 (seconds)	1000	100
Total time (seconds)	1001	110

شکل ۲-۵: زمان اجرای دو برنامه بر روی دو ماشین مختلف

طبق تعریفی که ما از سرعت داشتیم، عبارتهای زیر در مورد شکل ۲-۵ صادق هستند:

- ماشین A برای برنامه ۱، ۱۰ مرتبه سریعتر از ماشین B عمل می‌کند.
- ماشین B برای برنامه ۲، ۱۰ مرتبه سریعتر از ماشین A عمل می‌کند.

این دو جمله خیلی گیج کننده هستند. در نهایت کدام ماشین سریعتر است؟

ساده‌ترین راه برای خلاصه‌سازی کارآیی در این موارد که بیش از چند برنامه داریم این است که مجموع زمانهای اجرا را نسبت به هم مقایسه کنیم. بنابراین:

$$\frac{10 \times \text{مجموع زمان اجرای برنامه‌ها بر روی ماشین A}}{\text{کارآیی ماشین B}} = \frac{9 \times \text{مجموع زمان اجرای برنامه‌ها بر روی ماشین B}}{\text{کارآیی ماشین A}} = 9/1$$

یعنی اینکه ماشین B برای مجموعه برنامه‌های ۱ و ۲، ۹/۱ مرتبه سریعتر از ماشین A است.

این نوع خلاصه سازی مستقیماً متناسب با زمان اجرا می‌باشد. اگر workload ما به صورتی باشد که برنامه‌های ۱ و ۲ را به یک اندازه اجرا کنیم، رابطه فوق می‌تواند برای مقایسه ماشین‌های مختلف استفاده شود.

بعضی مواقع به جای رابطه فوق میانگین زمان اجرا که آن هم متناسب با زمان اجراست را به کار می‌برند. این میانگین، میانگین حسابی^{۱۴} یا AM نامیده می‌شود و از فرمول زیر بدست می‌آید:

$$AM = \frac{1}{n} \sum time_i$$

که در آن $time_i$ زمان اجرای برنامه i ام موجود در workload متشکل از n برنامه است. هر چقدر AM کوچکتر باشد، کارآیی ماشین بالاتر خواهد بود.

AM زمانی قابل استفاده است که برنامه‌های موجود در workload به یک اندازه اجرا شوند، اگر غیر از این باشد می‌توانیم برای هر برنامه یک وزنی اختصاص دهیم که نشان دهنده تعداد تکرار برنامه در آن workload باشد. به طور مثال اگر ۲۰ درصد کارها در workload برنامه ۱ و ۸۰ درصد کارها در workload برنامه ۲ باشد در این صورت وزنها به صورت ۰/۲ و ۰/۸ خواهند بود. با استفاده از جمع کردن حاصل ضرب فاکتورهای وزن در زمان اجرا فرمول میانگین حسابی وزن‌دار^{۱۵} یا WAM بدست می‌آید:

$$WAM = \sum_{i=1}^n weight_i \times time_i$$

¹⁴ - Arithmetic Mean

¹⁵ - Weighted Arithmetic Mean

فرمول میانگین حسابی حالت خاصی از فرمول میانگین حسابی وزن دار می باشد که در آن وزن همه برنامه ها مساوی است.

۲-۶- موضوع واقعی: benchmark های SPEC95 و کارایی پردازنده ها

SPEC معروفترین مجموعه benchmark های CPU می باشد. SPEC که مخفف عبارت System Performance Evaluation Cooperative (مجموعه ارزیابی کننده کارایی سیستم) می باشد، در سال ۱۹۸۹ توسط تعدادی از کمپانی های کامپیوتر برای این منظور ایجاد شد که اندازه گیری و خلاصه سازی کارایی از طریق یک پروسه اندازه گیری کنترل شده، بهبود داده شود و همچنین benchmark های واقعی تری برای اندازه گیری کارایی مورد استفاده قرار گیرند. تاکنون چهار نسخه از SPEC وارد بازار شده است: SPEC89، SPEC92، SPEC95 و SPEC2000. مجموعه SPEC95 نسخه سال ۹۵ SPEC می باشد که شامل ۸ برنامه اعداد صحیح و ۱۰ برنامه ممیز شناور می باشد. مجموعه برنامه های SPEC95 در شکل ۲-۶ نشان داده شده است.

Benchmark	Description
go	Artificial intelligence; plays the game of Go
m88ksim	Motorola 88k chip simulator; runs test program
gcc	The Gnu C compiler generating SPARC code
compress	Compresses and decompresses file in memory
li	Lisp interpreter
ijpeg	Graphic compression and decompression
perl	Manipulates strings and prime numbers in the special-purpose programming language Perl
vortex	A database program
tomcatv	A mesh generation program
swim	Shallow water model with 513 x 513 grid
su2cor	quantum physics; Monte Carlo simulation
hydro2d	Astrophysics; Hydrodynamic Navier Stokes equations
mgrid	Multigrid solver in 3-D potential field
applu	Parabolic/elliptic partial differential equations
trub3d	Simulates isotropic, homogeneous turbulence in a cube
apsi	Solves problems regarding temperature, wind velocity, and distribution of pollutant
fpppp	Quantum chemistry
wave5	Plasma physics; electromagnetic particle simulation

شکل ۲-۶: مجموعه benchmark های SPEC95

معمولاً دو گزارش جداگانه برای مجموعه برنامه های صحیح و برنامه های ممیز شناور داده می شود. روند اندازه گیری با استفاده از برنامه های مجموعه SPEC به این صورت است که در ابتدا عدد های اندازه گیری شده برای زمان اجرا با تقسیم کردن زمان اجرای برنامه بر روی ماشین SUN sparcc 10/40 به زمان اجرای بدست آمده بر روی ماشینی که می خواهیم کارایی آن را اندازه بگیریم، به صورت نرمال

شده درآورده می‌شود. سپس این عدد نرمال شده به یک معیار برای اندازه‌گیری کارایی منجر می‌شود که نسبت SPEC ratio نامیده می‌شود. هر چقدر عدد بدست آمده برای SPEC ratio بزرگتر باشد نشان دهنده کارایی بالاتر می‌باشد (در واقع SPEC ratio با زمان اجرا نسبت معکوس دارد). خلاصه- سازی‌های SPECint95 و SPECfp95 با استفاده از میانگین هندسی SPEC ratio ها بدست می‌آید که اولی مربوط به برنامه‌های صحیح و دومی مربوط به برنامه‌های ممیز شناور مجموعه SPEC می‌باشد. فرمول میانگین هندسی¹⁶ به صورت زیر است:

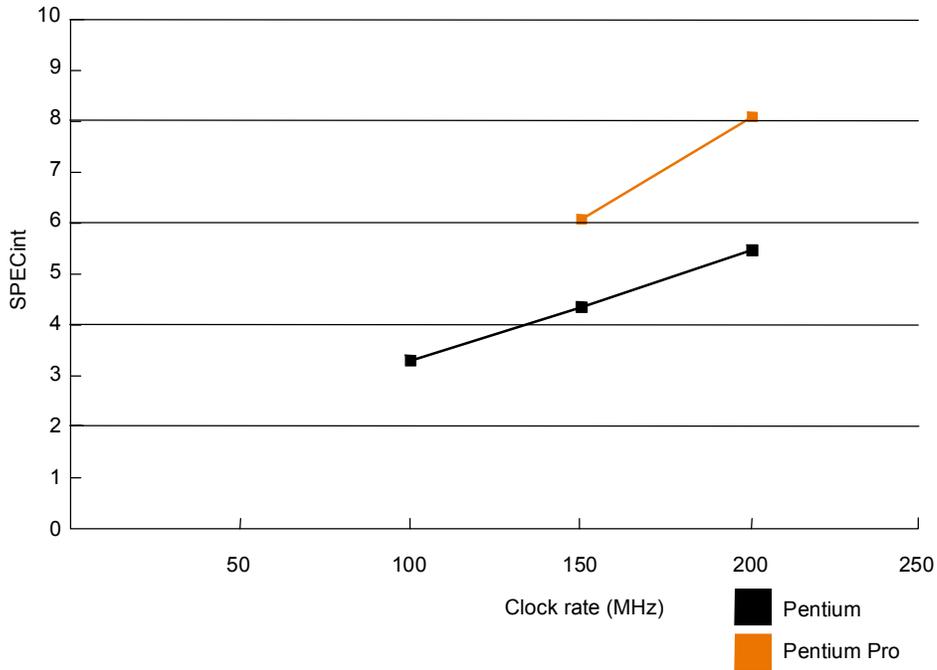
$$\sqrt[n]{\prod_{i=1}^n (execution\ time\ ratio)_i}$$

که در آن $(execution\ time\ ratio)_i$ زمان اجرای نرمال شده بر اساس یک ماشین مرجع برای برنامه i -ام مجموعه benchmark ها می‌باشد. همان طور که قبلاً نیز توضیح داده شده است برای یک ISA مفروض، افزایش کارایی CPU می‌تواند به یکی از سه روش زیر صورت بگیرد:

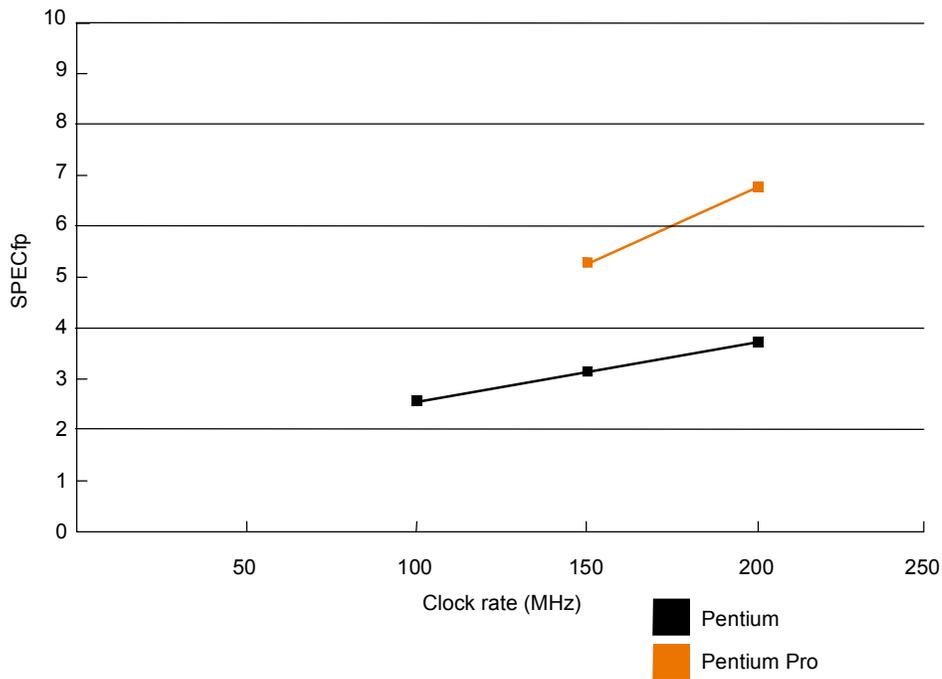
- ۱- افزایش فرکانس کلاک
- ۲- بهبود ساختار پردازنده که باعث کم شدن CPI شود
- ۳- بهبودهای کامپایلری که باعث کم شدن تعداد دستورات یا تولید دستوراتی با میانگین CPI پایین تر می‌شود

برای نشان دادن بهبودهای فوق بر روی کارایی، شکل‌های ۲-۷ و ۲-۸ به ترتیب اندازه‌گیری‌های SPECint95 و SPECfp95 را برای پردازنده‌های Pentium و Pentium pro نشان می‌دهند. چون SPEC نیاز به این دارد که حتماً روی سخت‌افزار واقعی اجرا شود و سیستم حافظه تأثیر زیادی بر روی کارایی دارد، سیستم‌های دیگری که بر اساس این پردازنده‌ها ساخته می‌شوند ممکن است مقداری با این شکلها متفاوت باشند چون سیستم حافظه آنها ممکن است قدری متفاوت باشد. ماشین‌های اینتل که در اینجا اندازه‌گیری شده‌اند سیستم‌های حافظه و کامپایلر قوی‌تری استفاده کرده‌اند و این بدان معناست که اکثر ماشین‌های دیگری که با استفاده از این پردازنده ساخته می‌شوند کارایی پایین‌تری برای benchmarks های SPEC خواهند داشت.

¹⁶ - Geometric Mean



شکل ۷-۲: SPECint محاسبه شده برای پردازنده‌های Pentium و Pentium pro در فرکانس‌های مختلف



شکل ۸-۲: SPECfp محاسبه شده برای پردازنده‌های Pentium و Pentium pro در فرکانس‌های مختلف

چند برداشت مهم از شکل‌های ۷-۲ و ۸-۲ قابل مشاهده است. مهمترین آنها بهبود کارایی پردازنده Pentium pro نسبت به Pentium می‌باشد. در یک فرکانس فرضی SPECint نشان می‌دهد که

Pentium pro، ۱/۴ تا ۱/۵ برابر سریعتر از Pentium می‌باشد و SPECfp95 نشان می‌دهد که Pentium pro، ۱/۷ تا ۱/۸ برابر سریعتر می‌باشد. اگرچه بهبودهای کامپایلری مشخص برای هر پردازنده‌ای وجود دارد، اکثر بهبودهای کارآیی برای Pentium pro از طریق بهبودهای ساختاری (سخت‌افزاری) ایجاد شده است.

برداشت مهم دیگری که از شکلها به دست می‌آید این است که وقتی که فرکانس کلاک با استفاده از ضریب مشخص اضافه می‌شود، کارآیی پردازنده با ضریب کمتری اضافه می‌شود. به طور مثال وقتی که فرکانس کلاک Pentium از ۱۰۰ مگاهرتز به ۲۰۰ مگاهرتز افزایش پیدا می‌کند (دو برابر می‌شود)، کارآیی SPECint 95 فقط ۱/۷ و کارآیی SPECfp95 فقط ۱/۴ برابر بهبود می‌یابد. دلیل این امر به سیستم حافظه پردازنده بر می‌گردد. در کل، چون سرعت حافظه اصلی تغییر چندانی نمی‌کند، اضافه شدن سرعت پردازنده در گلوگاه سیستم حافظه گرفتار می‌شود. این تأثیر در benchmark های ممیز شناور بیشتر مشهود است چون آنها نسبت به برنامه‌های integer بزرگتر هستند. این رفتار مثال خوبی برای قانون امدال می‌باشد که بیانگر میزان بهبود بر اساس میزان مؤثر بودن آن می‌باشد.

در مقام مقایسه، میزان افزایش کارآیی پردازنده Pentium pro، با اضافه شدن فرکانس کلاک بهتر از Pentium می‌باشد. هر چند که این میزان افزایش برابر ضریب افزایش کلاک نیست. به طور مثال، بهبود کارآیی SPECfp95 از فرکانس ۱۵۰ به ۲۰۰ مگاهرتز در Pentium pro، ۱/۲۴ می‌باشد در حالی که پردازنده Pentium بهبودی در حدود ۱/۱۸ را برای این بازه نشان می‌دهد. در فصل ۷ در مورد تأثیر سیستم حافظه در کارآیی پردازنده صحبت خواهد شد و اینکه چرا پردازنده Pentium pro سیستم حافظه بهتری دارد.

۲-۸ نکات پایانی

ما در این فصل بر روی performance یا همان کارآیی سیستم‌ها و نحوه اندازه‌گیری آن تمرکز کردیم. تمرکز کردن تنها بر روی کارآیی بدون در نظر گرفتن هزینه^{۱۷} خیلی منطقی نیست و ما باید بین این دو پارامتر تعادل برقرار کنیم.

¹⁷ - Cost

ما در این فصل دیدیم که با استفاده از زمان اجرای برنامه‌های واقعی بر روی پردازنده‌ها، یک روش قابل اطمینان برای اندازه‌گیری و گزارش کردن کارایی وجود دارد. این زمان اجرا وابسته به تعدادی فاکتور مهم می‌باشد که این وابستگی در فرمول زیر مشاهده می‌شود:

$$\text{پریود کلاک} \times \text{CPI} \times \text{تعداد دستورات} = \text{زمان اجرا}$$

ما به وفور از این فرمول و فاکتورهای آن بهره خواهیم گرفت. همان طور که توضیح داده شد هیچ کدام از این فاکتورها به تنهایی نمی‌توانند نشان دهنده کارایی باشند و فقط حاصل ضرب آنها که مساوی زمان اجراست قابل استفاده است. مطمئناً دانستن این فرمول به تنهایی برای طراحی کردن یا ارزیابی کردن یک کامپیوتر نمی‌تواند کافی باشد. ما باید بفهمیم که چگونه جنبه‌های مختلف طراحی، این فاکتورهای مهم را تحت تأثیر قرار می‌دهند. این نوع دانش شامل موارد متعددی می‌شود از جمله: تأثیر معماری مجموعه دستورات بر تعداد دستورات دینامیکی، تأثیر pipelining و سیستم حافظه بر روی CPI، و تأثیر تکنولوژی و سازمان یک کامپیوتر بر روی فرکانس کلاک. هنر طراحی کامپیوتر فقط اعمال کردن عدد در فرمول کارایی نیست بلکه مشخص کردن این موضوع هست که طراحی‌های مختلف چگونه کارایی و هزینه را تحت تأثیر قرار می‌دهند.

اکثر کاربران کامپیوتر به هر دو مورد هزینه و کارایی دقت می‌کنند. فهمیدن ارتباط بین جنبه‌های مختلف طراحی و کارایی آن کار مشکلی است و مشخص کردن هزینه قابلیت‌های مختلف یک طرح نیز کار به مراتب مشکل‌تری است. هزینه یک ماشین فقط به هزینه قسمت‌های تشکیل دهنده آن محدود نیست، بلکه هزینه‌هایی همچون هزینه‌های نیروی انسانی برای جمع‌آوری سیستم^{۱۸}، هزینه تحقیقات، هزینه توسعه سیستم، تبلیغات بازار و ... را نیز باید به آن اضافه نمود. امروزه به دلیل پیشرفت سریعی که در تکنولوژی‌های پیاده‌سازی ایجاد می‌شود، اهمیت دادن بیش از حد به مورد هزینه در طول مدت‌هایی حدود ۶ ماه یا یک سال خیلی منطقی به نظر نمی‌رسد.

و نکته آخر اینکه معماری‌های مختلف کامپیوتر توسط کارایی و هزینه سنجیده می‌شوند و پیدا کردن تعادل بین این دو همیشه به عنوان یک هنر در طراحی کامپیوتر مطرح بوده است.

مطالب مهم این فصل:

- Performance یا همان کارایی یک معیار مهم برای مقایسه سیستم‌هاست.
- دو روش معتبر برای اندازه‌گیری کارایی وجود دارد: محاسبه زمان اجرا و محاسبه throughput
- فرمول اصلی کارایی: پریود کلاک \times CPI \times تعداد دستورات = زمان اجرا

¹⁸ - Assemble

- بهترین راه برای اندازه‌گیری کارایی این است که از برنامه‌های واقعی استفاده کنیم.
- Benchmark ها برنامه‌هایی هستند که برای اندازه‌گیری کارایی انتخاب می‌شوند و باید به اندازه کافی بزرگ باشند که اجرای آنها توسط کامپایلر و یا سخت افزار بهینه سازی نگردد و تا حد ممکن از برنامه‌های واقعی انتخاب شوند. مهمترین مجموعه benchmark های شناخته شده مجموعه SPEC می‌باشد.
- قانون امدال بیان می‌کند که ما از یک بهبود داده شده در طرح چقدر می‌توانیم انتظار بهبود کارایی را داشته باشیم. میزان بهبود کارایی به میزان مؤثر بودن تغییر داده شده بستگی دارد.
- اگر به جای یک برنامه چند برنامه برای ارزیابی کارایی استفاده شود و بخواهیم کارایی را خلاصه سازی کنیم از فرمول‌های میانگین حسابی و یا میانگین حسابی وزن دار استفاده می‌کنیم. میانگین هندسی نیز قابل استفاده است ولی میانگین هندسی دارای ضعف‌هایی است که نمی‌توان آن را به عنوان یک معیار خوب قبول کرد.
- گزارش تهیه شده برای کارایی باید طوری باشد که توسط دیگران قابل تکرار باشد. بنابراین باید همه چیزهایی که برای تکرار یک آزمایش لازم است در گزارش آورده شود.